

The condensation transition in random hypergraph 2-coloring

Amin Coja-Oghlan* and Lenka Zdeborova†

September 27, 2011

Abstract

For many random constraint satisfaction problems such as random satisfiability or random graph or hypergraph coloring, the best current estimates of the threshold for the existence of solutions are based on the *first* and the *second moment method*. However, in most cases these techniques do not yield matching upper and lower bounds. Sophisticated but non-rigorous arguments from statistical mechanics have ascribed this discrepancy to the existence of a phase transition called *condensation* that occurs shortly before the actual threshold for the existence of solutions and that affects the combinatorial nature of the problem, rendering the second moment method powerless (Krzakala, Montanari, Ricci-Tersenghi, Semerjian, Zdeborova: PNAS 2007). In this paper we prove for the first time *that* a condensation transition exists in a natural random CSP, namely in random hypergraph 2-coloring. Perhaps surprisingly, we find that the second moment method breaks down strictly *before* the condensation transition. Our proof also yields slightly improved bounds on the threshold for random hypergraph 2-colorability. We expect that our techniques can be extended to other, related problems such as random k -SAT or random graph k -coloring.

Key words: random structures, phase transitions, hypergraph 2-coloring, second moment method.

*University of Warwick, Zeeman building, Coventry CV4 7AL, UK, a.coja-oghlan@warwick.ac.uk. Supported by EPSRC grant EP/G039070/2.

†Institut de Physique Théorique, IPhT, CEA Saclay, and URA 2306, CNRS, 91191 Gif-sur-Yvette, France, lenka.zdeborova@gmail.com.

1 Introduction and results

For many random constraint satisfaction problems such as random k -SAT, random graph coloring, or random hypergraph coloring the best current bounds on the thresholds for the existence of solutions derive from the *first* and the *second moment method*. However, in most cases these simple techniques do not yield matching upper and lower bounds. In effect, for most random CSPs the *precise* threshold for the existence of solutions remains unknown. Examples of this include random k -SAT, random graph k -coloring, and the 2-coloring problem in random k -uniform hypergraphs ($k \geq 3$).

In this paper we investigate the origin of this discrepancy with the example of the random hypergraph 2-coloring problem, in which the second moment analysis is technically relatively simple. First, we present an approach to improve slightly over the naive second moment argument. But more importantly, we establish the existence of a further phase transition below the threshold for the existence of solutions. At this so-called *condensation transition*, whose existence was predicted on grounds of sophisticated but non-rigorous statistical mechanics arguments [9, 16], the combinatorial nature of a ‘typical’ solution becomes significantly more complicated. Arguably, beyond the condensation transition it is conceptually more difficult to prove that solutions exist, and indeed in several random CSPs condensation seems to pose the key obstacle to determining the precise threshold for the existence of solutions. Here we prove rigorously for the first time that a condensation transition indeed exists.

To define the *random hypergraph 2-coloring* problem, let $V = \{1, \dots, n\}$ be a set of vertices, let $k \geq 3$, and let $H_k(n, m)$ be a random k -uniform hypergraph on V obtained by inserting a random set of m edges out of the $\binom{n}{k}$ possible edges. A 2-coloring of H is a map $\sigma : V \rightarrow \{0, 1\}$ such that no edge e of H is monochromatic. Throughout the paper, we will let $r = m/n$ denote the *density* of the random hypergraph. An event \mathcal{E} occurs *with high probability* (‘w.h.p.’) if its probability tends to one as $n \rightarrow \infty$. We let $\mathcal{S}(H)$ denote the set of all 2-colorings of the hypergraph H , and we let $Z(H) = |\mathcal{S}(H)|$.

Friedgut’s sharp threshold theorem implies that for any $k \geq 3$ there *exists* a threshold $r_{col} = r_{col}(k, n)$ such that for any $\varepsilon > 0$ the random hypergraph $H_k(n, m)$ of density $r = m/n < (1 - \varepsilon)r_{col}$ is 2-colorable w.h.p., while for $r > (1 + \varepsilon)r_{col}$ it is w.h.p. not [11, 12].¹ Although the precise threshold r_{col} is not known for any $k \geq 3$, the first and the second moment methods can be used to derive upper and lower bounds. To put our results in perspective, let us briefly recap these techniques.

The first and the second moment method. The first moment method yields an upper bound on r_{col} . More precisely, by Markov’s inequality,

$$\mathbb{P}[H_k(n, m) \text{ is 2-colorable}] = \mathbb{P}[Z \geq 1] \leq \mathbb{E}[Z].$$

Hence, if for some density r the first moment $\mathbb{E}[Z]$ satisfies $\mathbb{E}[Z] = o(1)$, then $r_{col} \leq r$ (for large enough n). Indeed, it is easy to compute $\mathbb{E}[Z]$ explicitly, and to verify that there is a critical density $r_{first} = 2^{k-1} \ln(2) - \ln(2)/2 + o_k(1)$ such that $\mathbb{E}[Z] = \exp(\Omega(n)) \gg 1$ if $r < r_{first}$ while $\mathbb{E}[Z] = o(1)$ if $r > r_{first}$. Hence, $r_{col} \leq r_{first}$.

Even though $\mathbb{E}[Z] = \exp(\Omega(n))$ is *exponentially* large in n for $r < r_{first}$, this does, of course, not imply that $H_k(n, m)$ is 2-colorable with high probability: it could simply be that a tiny number of hypergraphs drive up the expected number of 2-colorings because they possess excessively many of them. The purpose of the second moment method is to rule this possibility out. More precisely, the second moment argument is based on the *Paley-Zygmund inequality*

$$\mathbb{P}[H_k(n, m) \text{ is 2-colorable}] = \mathbb{P}[Z > 0] \geq \frac{\mathbb{E}[Z]^2}{\mathbb{E}[Z^2]}.$$

¹It is widely conjectured that $\lim_{n \rightarrow \infty} r_{col}(k, n)$ exists for any $k \geq 3$. Hence we will take the liberty of just speaking of ‘the threshold r_{col} ’ (for $k \geq 3$ given).

Hence, if for some density $r < r_{first}$ we can show that

$$\mathbb{E}[Z^2] \leq C \cdot \mathbb{E}[Z]^2 \quad (1)$$

with $C = C(k, r) > 0$ independent of n , then $\mathbb{P}[H_k(n, m) \text{ is 2-colorable}] \geq 1/C$. That is, the probability of 2-colorability is bounded away from 0 as $n \rightarrow \infty$. Therefore, the sharp threshold theorem implies that $r_{col} \geq r$. Indeed, Achlioptas and Moore [3] proved that there is a critical density $r_{second} = 2^{k-1} \ln 2 - (1 + \ln 2)/2 + o_k(1)$ such that (1) holds for all $r < r_{second}$ but is violated for $r > r_{second}$. In summary, the first/second moment arguments yield the bounds

$$r_{second} = 2^{k-1} \ln 2 - \frac{1 + \ln 2}{2} + o_k(1) \leq r_{col} \leq r_{first} = 2^{k-1} \ln 2 - \frac{\ln 2}{2} + o_k(1). \quad (2)$$

Approaching the condensation threshold. How could we improve the lower bound on r_{col} ? The second moment analysis in [3] is tight, and thus simply performing a better calculation will not suffice. Indeed, as observed in [3], for $r > r_{second}$ we have $\mathbb{E}[Z^2] > \exp(\Omega(n)) \cdot \mathbb{E}[Z]^2$, i.e., the second moment method fails *dramatically*. But why? One possibility could be that the *expectation* $\mathbb{E}[Z]$ is driven up by a tiny minority of hypergraphs with excessively many 2-colorings, i.e., that $Z \leq \exp(-\Omega(n))\mathbb{E}[Z]$ w.h.p. In this case (1) would fail to hold because the second moment $\mathbb{E}[Z^2]$ would exacerbate the contribution of the few ‘rich’ hypergraphs even more than the first moment. A second possibility is that Z is ‘close’ to $\mathbb{E}[Z]$ w.h.p., but without being sufficiently concentrated for (1) to hold. The following theorem, which improves the lower bound in (2) by an additive $(1 - \ln(2))/2 \approx 0.153$, shows that up to $r_{cond} = 2^{k-1} \ln 2 - \ln 2 > r_{second}$, the second scenario is true.

Theorem 1.1 *There is a constant $k_0 \geq 3$ such that for all $k \geq k_0$ and $r < r_{cond}$ the random hypergraph $H_k(n, m)$ is 2-colorable w.h.p. and*

$$\ln Z \sim \ln \mathbb{E}[Z] \quad \text{w.h.p.} \quad (3)$$

For $r < r_{cond}$ the *expected* number $\mathbb{E}[Z]$ of 2-colorings is exponentially large in n . Hence, (3) shows that for $r < r_{cond}$ w.h.p. Z is exponentially large as it coincides with $\mathbb{E}[Z]$ up to sub-exponential terms.

The proof of Theorem 1.1 is based on an enhanced second moment argument that takes the ‘geometry’ of the set $\mathcal{S}(H_k(n, m))$ of 2-colorings of the random hypergraph into account. As a corollary of this argument, we obtain a result on the ‘shape’ of this set, viewed as a subset of the n -dimensional Hamming cube $\{0, 1\}^n$ equipped with the Hamming distance. To state this result, let us say that a 2-coloring σ of a hypergraph H on n vertices is (α, β, γ) -*shattered* for $\alpha, \gamma > 0$ and $\beta > \alpha$ if the following is true.

SH1. There is no 2-coloring $\tau \in \mathcal{S}(H)$ with $\alpha n < \text{dist}(\sigma, \tau) < \beta n$.

SH2. The set $\mathcal{C}_\alpha(\sigma)$ of all 2-colorings $\tau \in \mathcal{S}(H)$ with $\text{dist}(\sigma, \tau) \leq \alpha n$ has size $|\mathcal{C}_\alpha(\sigma)| \leq \exp(-\gamma n)Z(H)$.

Intuitively, this means that σ is part of a ‘cluster’ $\mathcal{C}_\alpha(\sigma)$ of 2-colorings, whose size is exponentially smaller than the total number $Z(H)$ of 2-colorings. Furthermore, there is a ‘gap’ of size $(\beta - \alpha)n$ between this cluster and the remaining 2-colorings of H .

Corollary 1.2 *There is a constant $k_0 \geq 3$ such that for any $k \geq k_0$ there is $\gamma_k > 0$ such that for $r < r_{cond}$ all 2-colorings of the random hypergraph $H_k(n, m)$ are $(0.01, 0.49, \gamma_k)$ -shattered w.h.p.*

Corollary 1.2 implies that w.h.p. the set of 2-colorings of $H = H_k(n, m)$ has a decomposition $\mathcal{S}(H) = \bigcup_{i=1}^N S_i$ into subsets that each comprise only an exponentially small fraction of all 2-colorings and that are mutually at Hamming distance at least $0.48n$. (Indeed, inductively choose S_i to be the local cluster $\mathcal{C}_{0.01}(\sigma)$)

of some 2-coloring $\sigma \notin \bigcup_{j < i} S_j$.) This decomposition allows us to explain intuitively why the ‘vanilla’ second moment argument fails for $r_{\text{second}} < r < r_{\text{cond}}$. In fact, we can write $Z(H)^2 = \sum_{i,j=1}^N |S_i| \cdot |S_j|$. To estimate the expectation of this quantity, we need to bound on the number N of components and their sizes $|S_i|$. As we will see in Section 4, the naive second moment argument overestimates the ‘cluster sizes’ $|S_i|$ grossly. We overcome this problem by investigating the internal structure of the ‘clusters’ S_i . We expect that this approach extends to other problems such as random k -SAT or random graph k -coloring, although the technical details will be far more intricate.

Into the condensation phase. As we will see next, even the enhanced second moment argument from Theorem 1.1 does not give the precise threshold for 2-colorability. The intuitive reason is that for densities beyond r_{cond} , the *expected* number $\mathbb{E}[Z]$ of 2-colorings is indeed driven up excessively by a tiny minority of hypergraphs with an abundance of 2-colorings.

Theorem 1.3 *There exist a constant $k_0 \geq 3$ and a sequence $\varepsilon_k \rightarrow 0$ such that for any $k \geq k_0$ there are $\delta_k > 0, \zeta_k > 0$ such that the following two statements are true.*

1. *W.h.p. $H_k(n, m)$ is 2-colorable for all $r < r_{\text{cond}} + \varepsilon_k + \delta_k$.*
2. *For any density r with $r_{\text{cond}} + \varepsilon_k < r < r_{\text{col}}$ we have*

$$\ln Z < \ln \mathbb{E}[Z] - \zeta_k n \quad \text{w.h.p.} \quad (4)$$

The second statement asserts that for densities between $r_{\text{cond}} + \varepsilon_k$ and the actual (unknown) 2-colorability threshold r_{col} , the *expected* number $\mathbb{E}[Z]$ of 2-colorings exceeds the *actual* number Z by an exponential factor $\exp(\zeta_k n)$ w.h.p. This contrasts with Theorem 1.1, which shows that below r_{cond} , Z is of the same exponential order as $\mathbb{E}[Z]$ w.h.p. Furthermore, the first part of Theorem 1.3 ensures that the regime of densities where (4) holds is non-empty, as the true threshold r_{col} is indeed strictly greater than $r_{\text{cond}} + \varepsilon_k$. This so-called *condensation transition* at density $r_{\text{cond}} = 2^{k-1} \ln 2 - \ln 2$ was predicted on the basis of non-rigorous statistical mechanics arguments [9, 16].

In mathematical physics, the term ‘phase transition’ is usually defined as a point where the function $F(r) = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E}[\ln(1 + Z)]$ is non-analytic. However, it is not currently known if the limit $F(r)$ exists. (Bayati, Gamarnik and Tetali [7] proved that for any density r , the corresponding limit of the partition function at any fixed positive temperature exists.) It is not difficult to see that Theorems 1.1 and 1.3 imply that around $r = r_{\text{cond}}$, the function $F(r)$ in fact is non-analytic if the limit exists (because for $r < r_{\text{cond}}$, $F(r)$ coincides with the *linear* function $\lim_{n \rightarrow \infty} \frac{1}{n} \ln \mathbb{E}[Z]$).

The term ‘condensation’ is meant to express that w.h.p. the set $\mathcal{S}(H_k(n, m))$ of all 2-colorings has a drastically different shape than in the ‘shattered’ regime of Corollary 1.2. To express this, let us call a 2-coloring of a hypergraph H on n vertices (α, β, γ) -*condensed* if

CO1. There is no 2-coloring $\tau \in \mathcal{S}(H)$ with $\alpha n < \text{dist}(\sigma, \tau) < \beta n$.

CO2. The set $\mathcal{C}_\alpha(\sigma)$ of all 2-colorings $\tau \in \mathcal{S}(H)$ with $\text{dist}(\sigma, \tau) \leq \alpha n$ has size $|\mathcal{C}_\alpha(\sigma)| \geq \exp(-\gamma n) Z(H)$.

(The difference between **SH1–SH2** and the above is that **CO2** imposes a *lower* bound on $|\mathcal{C}_\alpha(\sigma)|$.)

Corollary 1.4 *There exist a constant $k_0 \geq 3$ and a sequence $\varepsilon_k \rightarrow 0$ such that for any $k \geq k_0$ there exist a sequence $r(n)$ of densities satisfying $|r(n) - r_{\text{cond}}| \leq \varepsilon_k$ such that $H_k(n, m)$ with $m = r(n) \cdot n$ has the following two properties w.h.p.*

1. $H_k(n, m)$ is 2-colorable.
2. A random 2-coloring $\sigma \in \mathcal{S}(H_k(n, m))$ is $(0.01, 0.49, o(1))$ -condensed w.h.p.

This means that at a particular density $r(n)$, i.e., *right at* the condensation transition, the size of the local cluster of a ‘typical’ 2-coloring σ of $H_k(n, m)$ satisfies $\ln |\mathcal{C}_{0.01}(\sigma)| \sim \ln Z$ w.h.p. In other words, the size of the cluster of a ‘typical’ 2-coloring has the same exponential order as the set of *all* 2-colorings. This contrasts with the ‘shattered’ scenario of Corollary 1.2, where w.h.p. *all* clusters only comprise an exponentially small fraction of the entire set $\mathcal{S}(H_k(n, m))$. The statistical physics work [9, 16] suggests that indeed, the conclusions of Corollary 1.4 hold in the *entire* regime between the condensation transition and the 2-colorability threshold.

Discussion. The significance of the slightly better lower bound on the threshold for hypergraph 2-colorability provided by Theorem 1.1 is that it allows us to prove the existence of the *condensation transition*. Beyond the condensation transition, the combinatorial nature of the problem becomes far more complicated. To see why, consider the following random experiment with $r < r_{col}$ (so that $H_k(n, m)$ is 2-colorable w.h.p.).

G1. Choose a random hypergraph $H = H_k(n, m)$, conditional on H being 2-colorable.

G2. Choose a 2-coloring $\sigma \in \mathcal{S}(H)$ uniformly at random and output (H, σ) .

The above experiment induces a probability distribution $g_{k,n,m}$ on the set $\Lambda_k(n, m)$ of hypergraph/2-coloring pairs that we call the *Gibbs distribution*. For $r < r_{col}$ the experiment corresponds to sampling a random 2-coloring of a random hypergraph, and thus understanding the above experiment is the key to studying the combinatorial nature of the hypergraph 2-colorability problem. But the experiment seems genuinely difficult to analyze. In fact, even for densities $r = O(2^{k-1}/k)$ far below the threshold for 2-colorability, it is not known how to efficiently construct, let alone sample, a 2-coloring of a random hypergraph [2].

But there is a related experiment called the *planted model* that is rather easy to implement and to study.

P1. Choose $\sigma \in \{0, 1\}^n$ uniformly at random.

P2. Choose a hypergraph $H = H_k(n, m, \sigma)$ with m edges uniformly at random among all hypergraphs for which σ is a proper 2-coloring, and output (H, σ) .

Let $p_{k,n,m}$ denote the distribution on $\Lambda_k(n, m)$ induced by **P1–P2**. It is not difficult to show that prior to the condensation phase, the distributions induced by the two experiments are ‘close’.

Proposition 1.5 ([1]) *Suppose that $r < r_{first}$ is such that $\ln Z \sim \ln \mathbb{E}[Z]$ w.h.p. Then*

$$\ln(g_{k,n,m}[\mathcal{B} | \{\ln Z \sim \ln \mathbb{E}[Z]\}]) \leq \ln(p_{k,n,m}[\mathcal{B}]) + o(n) \quad \text{for any event } \mathcal{B} \neq \emptyset. \quad (5)$$

The relationship (5) allows us to bound the probability of some ‘bad’ event \mathcal{B} in the Gibbs distribution by estimating its probability in the planted distribution. Indeed, Proposition 1.5 was used in [1] to study various properties of ‘typical’ 2-colorings of $H_k(n, m)$. In combination with Theorem 1.1 and the methods of [1], Proposition 1.5 can be used to get a pretty good idea what a 2-coloring of the random hypergraph $H_k(n, m)$ “typically looks like” *before* the condensation transition.

But beyond the condensation transition, all bets are off. As Theorem 1.3 shows, in the condensed regime we have $\ln Z < \ln \mathbb{E}[Z] - \Omega(n)$ w.h.p., i.e., the assumption of Proposition 1.5 is violated. Roughly speaking, the gap $\ln Z < \ln \mathbb{E}[Z] - \Omega(n)$ implies that a pair chosen from the planted distribution **P1–P2** corresponds to a pair chosen from the Gibbs distribution only with exponentially small probability. In fact, for densities beyond the condensation transition our proof of Theorem 1.3 exhibits an event \mathcal{B} for which (5) is violated, i.e., the planted model is no longer a good approximation to the Gibbs distribution. Furthermore, the statistical mechanics cavity technique suggests that getting a handle on the Gibbs measure (or other related measures) is far more complicated in the condensation phase. Overcoming this obstacle appears to be the remaining challenge to obtaining the precise threshold for hypergraph 2-colorability. The statistical mechanics reasoning [9, 16] suggests

Conjecture 1.6 *There is $\varepsilon_k \rightarrow 0$ such that $r_{col} \sim 2^{k-1} \ln 2 - (\frac{\ln 2}{2} + \frac{1}{4}) + \varepsilon_k$.*

One limitation of our approach is that we need to assume that $k \geq k_0$ is sufficiently big (whereas the standard second moment argument [3] applies to any $k \geq 3$). We need the lower bound on k to carry out a sufficiently accurate analysis of combinatorial structure of the solution space $\mathcal{S}(H_k(n, m))$. No attempt has been made to compute (let alone optimize) k_0 or the various other constants.

2 Related work

The two inequalities in (2) state the best previous bounds on the threshold for hypergraph 2-colorability from the paper of Achlioptas and Moore [3], which provided the prototype for the second moment analyses in other sparse random CSPs (e.g., [4, 5]). Since the second moment method is non-constructive, there is the separate algorithmic question: for what densities can a 2-coloring of a random hypergraph be constructed in polynomial time w.h.p.? The best current algorithm is known to succeed up to $r = c \cdot 2^{k-1}/k$ for some constant $c > 0$, i.e., up to a factor of about k below the 2-colorability threshold [2].

In [1] the geometry of the set $\mathcal{S}(H_k(n, m))$ of 2-colorings of the random hypergraph was investigated (among other things). It was shown that $\mathcal{S}(H_k(n, m))$ shatters into exponentially small well-separated ‘clusters’ for densities $(1 + \varepsilon_k)2^{k-1} \ln(k)/k < r < r_{second}$. Corollary 1.2 extends this picture up to $r < r_{cond}$. In addition, [1] also proved that in the regime $(1 + \varepsilon_k)2^{k-1} \ln(k)/k < r < r_{second}$ a typical 2-coloring σ of $H_k(n, m)$ is *rigid* w.h.p. in the sense that for most vertices v any 2-coloring τ with $\sigma(v) \neq \tau(v)$ has Hamming distance $\Omega(n)$ from σ . Our analysis, most notably the study of the structure of a typical ‘local cluster’ in Section 5, builds substantially on the concepts of shattering and rigidity from [1], but we will have to elaborate them in considerably more detail to get close quantitative estimates.

In many random CSPs other than random hypergraph 2-coloring the best current bounds on the thresholds for the existence of solutions derive from the second moment method as well. The most prominent examples are random graph k -coloring [4] and random k -SAT [5]. But the second moment argument extends naturally to a range of ‘symmetric’ random CSPs [17]. It would be interesting to see if/how our techniques can be generalized in order to prove the existence of a condensation transition in these other problems, particularly random graph k -coloring. However, since even the standard second moment analysis is quite involved in this case of random graph k -coloring, such a generalization will be technically challenging.

The random k -SAT problem is conceptually different because it is not ‘symmetric’. More precisely, in random hypergraph 2-coloring the *inverse* $1 - \sigma$ of a 2-coloring σ is a 2-coloring as well. This symmetry, which greatly simplifies the second moment argument, is absent in random k -SAT. As a consequence, as elaborated in [3, 5], in k -SAT the bound $\mathbb{E}[Z^2] = O(\mathbb{E}[Z]^2)$ does not hold for *any* density. Roughly speaking, to overcome this problem [5] focuses on a special type of satisfying assignments (“balanced” ones), whose number Z_* satisfies $\mathbb{E}[Z_*^2] = O(\mathbb{E}[Z_*]^2)$. Technically, this is accomplished by weighting satisfying assignments cleverly. While our techniques can be extended easily to establish the existence of a condensation transition for these *balanced* satisfying assignments in random k -SAT, this does not imply that condensation occurs with respect to the bigger set of *all* satisfying assignments. This would require a new approach for the direct analysis of the *total* number of satisfying assignments in random k -SAT.

We emphasize that our techniques are quite different from the ‘weighted’ second moment method in [5]. Indeed, the ‘asymmetry’ that motivated the weighting scheme in [5] is absent in random hypergraph 2-coloring. Instead of weighting, we employ a new idea that exploits the combinatorial structure of the ‘clusters’ into which the set $\mathcal{S}(H_k(n, m))$ of 2-colorings decomposes.

An example of a random CSP in which the precise threshold for the existence of solutions is known is random k -XORSAT. In this problem a second moment argument yields the precise thresholds (after ‘pruning’ the underlying hypergraph) [10, 19]. The explanation for this success is that random k -XORSAT does

not have a condensation phase due to the algebraic nature of the problem. Similarly, in random k -SAT with $k > \log_2 n$ (i.e., the clause length is *growing* with n) there is no condensation phase and, in effect, the second moment method yields the precise satisfiability threshold [8, 13]. A further class of problems where the condensed phase is conjectured to be empty are the ‘locked’ problems of [22].

In statistical mechanics the condensation transition was first predicted (using non-rigorous techniques) for the random k -SAT and the random graph k -coloring problems [16]. For random hypergraph 2-coloring the statistical mechanics prediction for the condensation threshold was derived in [9]. The structure of the condensed phase is described using a non-rigorous framework called *one-step replica symmetry breaking*. Interestingly, it was also conjectured that the structure of the condensed phase for large k is very similar to the structure of the random subcube model [18]. Our proofs verify this for random hypergraph 2-coloring.

Random CSPs, including random hypergraph 2-coloring, have been studied in statistical mechanics as models of disordered systems (such as glasses) under the name ‘diluted mean field models’. In this context the condensation transition corresponds to the so-called *Kauzmann transition* [15]. The present paper provides the first rigorous proof that this phase transition actually exists in a ‘diluted mean field model’.

3 Preliminaries

We need the following Chernoff bound on the tails of a binomially distributed random variable from [14, p. 21]. Let $\varphi(x) = (1+x) \ln(1+x) - x$. [14, p. 26]

Lemma 3.1 *Let X be a binomial random variable with mean $\mu > 0$. Then for any $t > 0$ we have*

$$\begin{aligned} \mathbb{P}[X > \mathbb{E}[X] + t] &\leq \exp(-\mu \cdot \varphi(t/\mu)), \\ \mathbb{P}[X < \mathbb{E}[X] - t] &\leq \exp(-\mu \cdot \varphi(-t/\mu)). \end{aligned}$$

In particular, for any $t > 1$ we have $\mathbb{P}[X > t\mu] \leq \exp[-t\mu \ln(t/e)]$.

The following large deviations principle for the binomial distribution can be found, e.g., in [14, p. 27].

Lemma 3.2 *Let $X = \text{Bin}(n, p)$ be a binomial random variable with $\mu = np > 0$. Let t be such that $\mu + t \in \{1, \dots, n-1\}$. Then*

$$\ln \mathbb{P}[X = \mu + t] \sim -\mu \varphi(t/\mu) - (n - \mu) \varphi(t/(n - \mu)).$$

The following is a mild generalization of ‘Laplace lemmas’ statements in [3, 10].

Lemma 3.3 *Let $\psi \in C^3(0, 1)$ be such that $\lim_{x \rightarrow 0} \psi(x) = \lim_{x \rightarrow 1} \psi(x) = 0$. Assume that $z \in (0, 1)$ is the unique global maximum of ψ , that $\psi(z) > 0$, and that $\psi''(z) < 0$. Then*

$$\sum_{d=1}^{n-1} \exp(\psi(d/n)) \leq O(\sqrt{n}) \exp(n\psi(z)).$$

Proof. Since $\psi \in C^3(0, 1)$, Taylor’s formula shows that

$$\psi(z + \delta) - \psi(z) = \frac{\delta^2}{2} \cdot \psi''(z) + O_\delta(\delta^2). \quad (6)$$

Moreover, as $\lim_{x \rightarrow 0} \psi(x) = \lim_{x \rightarrow 1} \psi(x) = 1$, for any fixed $\delta > 0$ we have

$$\sum_{d=1}^{n-1} \exp(\psi(d/n)) \sim \sum_{(z-\delta)n < d < (z+\delta)n} \exp(\psi(d/n)). \quad (7)$$

Suppose that $(z - \delta)n < d < (z + \delta)n$. Then (6) implies that for small enough δ ,

$$\begin{aligned} \exp[n\psi(d/n)] &= \exp[n\psi(z)] \cdot \exp \left[n \left(\frac{\psi''(z)}{2} \left(\frac{d}{n} - z \right)^2 + O((d/n - z)^3) \right) \right] \\ &\leq \exp[n\psi(z)] \cdot \exp \left[\frac{\psi''(z)}{3} \cdot \frac{(d - zn)^2}{n} \right]. \end{aligned} \quad (8)$$

Combining (7) and (8) yields the assertion. \square

The following lemma is implicit in [1].

Lemma 3.4 *For any $\varepsilon > 0$ and any $k \geq 3$ the following is true. Suppose that $r < r_{\text{cond}}$. Then w.h.p. $H_k(n, m)$ is such that*

$$\ln Z \sim \mathbb{E} \ln [1 + Z].$$

4 The enhanced second moment argument

In the rest of this paper, we assume that $k \geq k_0$ for some large enough constant k_0 . Moreover, to avoid floor and ceiling signs, we assume that n is even.

4.1 The local cluster and the demise of the vanilla second moment argument

We begin by briefly reviewing the ‘vanilla’ second moment method from [3]. This will provide the background for the enhanced the second moment argument that yields Theorem 1.1. As a first step, we need to work out the *expected* number $\mathbb{E}[Z]$ of 2-colorings.

Lemma 4.1 *We have $\mathbb{E}[Z] \sim 2^n(1 - 2^{1-k})^m$.*

Proof. Any fixed $\sigma : V \rightarrow \{0, 1\}$ is a 2-coloring of $H_k(n, m)$ iff $H_k(n, m)$ does not feature an edge that consists of vertices in one color class $\sigma^{-1}(i)$ only ($i = 1, 2$). In other words, σ ‘forbids’ $\binom{|\sigma^{-1}(0)|}{k} + \binom{|\sigma^{-1}(1)|}{k}$ out of the $\binom{n}{k}$ possible edges. Clearly, the number of ‘forbidden’ edges is minimized if both color classes $\sigma^{-1}(0), \sigma^{-1}(1)$ are the same size $n/2$. Furthermore, for all but a $o(1)$ -fraction of all 2^n possible $\sigma : V \rightarrow \{0, 1\}$ it is indeed true that both color classes have size $(1 \pm o(1))n/2$. By the linearity of expectation,

$$\mathbb{E}[Z] \sim 2^n \left(\binom{n}{k} - 2 \binom{n/2}{k} \right) / \binom{n}{m} \sim 2^n (1 - 2^{1-k})^m,$$

where the last step follows from Stirling’s formula. \square

Our goal is to identify the regime of densities r where $\mathbb{E}[Z^2] = O(\mathbb{E}[Z]^2)$, i.e., where the second moment method ‘works’. A technical issue is that Z includes 2-colorings σ whose color classes have (very) different sizes. To simplify our calculations we are going to confine ourselves to colorings σ whose color classes $\sigma^{-1}(0), \sigma^{-1}(1)$ have the same size. More precisely, let us call $\sigma : V \rightarrow \{0, 1\}$ *equitable* if

$|\sigma(0)| = |\sigma(1)| = n/2$, and let Z_e be the number of equitable 2-colorings of $H_k(n, m)$. Using Stirling's formula and, once more, the linearity of the expectation, it is not difficult to compute $E[Z_e]$: we have

$$E[Z_e] \sim \sqrt{\frac{2}{\pi}} \cdot \frac{2^n}{\sqrt{n}} (1 - 2^{1-k})^m = \Theta(1/\sqrt{n}) \cdot E[Z]. \quad (9)$$

Now, for what r do we have $E[Z_e^2] = O(E[Z_e]^2)$? We use the following elementary relation.

Fact 4.2 *For any equitable $\sigma : V \rightarrow \{0, 1\}$ we have $E[Z_e^2] = E[Z_e] \cdot E[Z_e | \sigma \text{ is a 2-coloring}]$.*

Proof. As $E[Z_e^2]$ equals the expected number of *pairs* of equitable 2-colorings, we find

$$\begin{aligned} E[Z_e^2] &= \sum_{\sigma, \tau} P[\sigma \text{ is a 2-coloring}] \cdot P[\tau \text{ is a 2-coloring} | \sigma \text{ is a valid 2-coloring}] \\ &= \sum_{\sigma} P[\sigma \text{ is a 2-coloring}] \cdot E[Z_e | \sigma \text{ is a 2-coloring}]. \end{aligned}$$

By symmetry, $E[Z_e | \sigma \text{ is a 2-coloring}]$ is the same for *all* equitable σ . Moreover, by the linearity of the expectation we have $E[Z_e] = \sum_{\sigma} P[\sigma \text{ is a 2-coloring}]$. \square

Thus, we need to compute $E[Z_e | \sigma \text{ is a 2-coloring}]$. In other words, for a *fixed* equitable $\sigma \in \{0, 1\}^n$ we need to study the random hypergraph $H_k(n, m)$ *given that* σ is a 2-coloring. This conditional distribution can be expressed easily: just choose a set of m edges uniformly at random from all edges that are bichromatic under σ (cf. step **P2** of the ‘planted model’ above). Let $H_k(n, m, \sigma)$ denote the resulting random hypergraph. Furthermore, given σ , let $Z_e(d)$ be the number of equitable 2-colorings τ with Hamming distance $\text{dist}(\sigma, \tau) = d$. Similarly, let $Z(d)$ be the *total* number of 2-colorings τ with $\text{dist}(\sigma, \tau) = d$. Then

$$E[Z_e | \sigma \text{ is a 2-coloring}] = \sum_{d=0}^n E_{H_k(n, m, \sigma)}[Z_e(d)] \leq \sum_{d=0}^n E_{H_k(n, m, \sigma)}[Z(d)]. \quad (10)$$

Fact 4.3 ([3]) *For any $0 < d < n$ we have*

$$\begin{aligned} E_{H_k(n, m, \sigma)}[Z(d)] &= \Theta(\sqrt{n/(d \cdot (n-d))}) \cdot \exp(\psi(d/n)), \\ E_{H_k(n, m, \sigma)}[Z_e(d)] &= \Theta(n/(d \cdot (n-d))) \cdot \exp(\psi(d/n)), \quad \text{with} \\ \psi &= \psi_{k,r} : (0, 1) \rightarrow \mathbf{R}, \quad x \mapsto -x \ln(x) - (1-x) \ln(1-x) + r \cdot \ln \left[1 - \frac{1-x^k - (1-x)^k}{2^{k-1} - 1} \right]. \end{aligned}$$

Fact 4.3 and (10) reduce the problem of computing $E[Z_e | \sigma \text{ is a 2-coloring}]$ (and thus $E[Z_e^2]$) to an exercise in calculus: we just need to study the function ψ .

Lemma 4.4 ([3]) *Suppose $r < r_{\text{first}}$. The function ψ satisfies $\psi(1/2) \sim \frac{1}{n} \ln E[Z]$, $\psi(1-x) = \psi(x)$, $\psi'(1/2) = 0$, and $\psi''(1/2) < 0$. Moreover,*

1. *if $\psi(1/2) > \psi(x)$ for all $x \in (0, 1)$, then $E[Z_e | \sigma \text{ is a 2-coloring}] \leq O(E[Z_e])$.*
2. *if there is some $x \in (0, 1)$ with $\psi(x) > \psi(1/2)$, then $E[Z_e | \sigma \text{ is a 2-coloring}] > E[Z] \cdot \exp(\Omega(n))$.*

Lemma 4.4 shows that the second moment method ‘works’ if and only if r is such that the function ψ takes its global maximum at $\frac{1}{2}$. Thus, let r_{second} be the supremum of all $r > 0$ with this property. Using basic calculus, one verifies that $r_{\text{second}} = 2^{k-1} \ln 2 - \frac{1}{2}(1 + \ln 2) + o_k(1)$ (see [3, Section 7]), and that for $r > r_{\text{second}}$ the function ψ attains its maximum, strictly greater than $\psi(1/2)$, in the interval $(0, 2^{-k/2})$. In effect, the second part of Lemma 4.4 shows that $E[Z^2] \geq E[Z_e^2] \geq \exp(\Omega(n))E[Z]$ for $r > r_{\text{second}}$, i.e., the ‘vanilla’ second moment argument breaks beyond r_{second} .

4.2 Improving the second moment argument: proof of Theorem 1.1

To improve over the naive second moment argument, we take another look at the function ψ . Let $\alpha = 2^{-k/2}$. Once more using basic calculus (see Section 4.4), we find

Lemma 4.5 *Suppose that $r_{\text{second}} < r < r_{\text{first}}$.*

1. *We have $\sup_{0 < x < \alpha} \psi(x) > \psi(1/2) > 0$.*
2. *For all $x \in (\alpha, 1/2 - \alpha) \cup (1/2 + \alpha, 1 - \alpha)$ we have $\psi(x) < -\psi(1/2) < 0$.*
3. *In the interval $[\alpha, 1 - \alpha]$ the function ψ attains its unique maximum at $1/2$.*

Lemma 4.5 allows us to deduce important information on the geometry of the set $\mathcal{S}(H_k(n, m, \sigma))$ of 2-colorings (similar arguments as the following have been used in [1] to prove that the set of all 2-colorings of $H_k(n, m)$ shatters into exponentially many well-separated pieces for a certain r). Indeed, combining Fact 4.3 and Lemma 4.5, we see that for distances $\alpha n < d \leq (\frac{1}{2} - \alpha)n$, the *expected* number of 2-colorings at distance d from σ is exponentially small:

$$\mathbb{E}_{H_k(n, m, \sigma)} [Z(d)] = \exp((1 + o(1))\psi(d/n)n) \leq \exp(-\Omega(n)).$$

Hence, $H_k(n, m, \sigma)$ does not have *any* 2-coloring τ such that $\text{dist}(\sigma, \tau) \in (\alpha n, (\frac{1}{2} - \alpha)n)$ w.h.p. Similarly, w.h.p. there is no 2-coloring τ with $\text{dist}(\sigma, \tau) \in ((\frac{1}{2} + \alpha)n, (1 - \alpha)n)$. Thus, w.h.p. the set of 2-colorings of $H_k(n, m, \sigma)$ decomposes into the ‘local cluster’

$$\mathcal{C}(\sigma) = \{\tau \in \mathcal{S}(H_k(n, m, \sigma)) : \text{dist}(\sigma, \tau) \leq \alpha n\}$$

of colorings ‘close’ to σ , the corresponding inverse colorings $\{\mathbf{1} - \tau : \tau \in \mathcal{C}(\sigma)\}$, and the remaining colorings τ with $\frac{1}{2} - \alpha \leq \text{dist}(\sigma, \tau)/n \leq \frac{1}{2} + \alpha$.

With this picture in mind, we can interpret the maximum of ψ in $(0, \alpha)$ as the *expected* size of the local cluster. More precisely, by Fact 4.3,

$$\mathbb{E}_{H_k(n, m, \sigma)} |\mathcal{C}(\sigma)| = \sum_{0 \leq d \leq \alpha n} \mathbb{E}_{H_k(n, m, \sigma)} [Z_d(\sigma)] = \exp \left[(1 + o(1))n \cdot \sup_{0 < x < \alpha} \psi(x) \right]. \quad (11)$$

Hence, the ‘vanilla’ second moment argument breaks down for $r > r_{\text{second}}$ because the *expected* size of the local cluster in $H_k(n, m, \sigma)$ exceeds the expected number $\mathbb{E}[Z]$ of 2-colorings in $H_k(n, m)$.

Our improvement over the plain second moment argument rests on the observation that for densities $r > r_{\text{second}}$ the *expected* size $\mathbb{E}|\mathcal{C}(\sigma)|$ exaggerates the *typical* size of the local cluster. More precisely, in Section 5 below we will investigate the combinatorial structure of the ‘planted’ formula $H_k(n, m, \sigma)$ closely to prove the following key fact.

Proposition 4.6 *Let $\sigma \in \{0, 1\}^V$ be equitable. If $r < r_{\text{cond}}$, then w.h.p. in the random formula $H_k(n, m, \sigma)$ the set $\mathcal{C}(\sigma) = \{\tau \in \mathcal{S}(H_k(n, m, \sigma)) : \text{dist}(\sigma, \tau) \leq 2^{-k/2}n\}$ has size $|\mathcal{C}(\sigma)| \leq \mathbb{E}[Z_e]$.*

Fix a density $r < r_{\text{cond}}$. Let us call a 2-coloring σ of a hypergraph H *good* if σ is equitable and the its local cluster $\mathcal{C}(\sigma) = \{\tau \in \mathcal{S}(H) : \text{dist}(\sigma, \tau) \leq 2^{-k/2}n\}$ has size $|\mathcal{C}(\sigma)| \leq \mathbb{E}[Z_e]$. Furthermore, let Z_g be the number of good 2-colorings of $H_k(n, m)$.

Corollary 4.7 *For any $r < r_{\text{cond}}$ we have $\mathbb{E}[Z_g] \sim \mathbb{E}[Z_e] = \Theta(n^{-1/2})\mathbb{E}[Z]$.*

Proof. Let \mathcal{H} be the set of all k -uniform hypergraphs on $V = \{1, \dots, n\}$ with precisely m edges. Let Λ_e be the set of all pairs (H, σ) with $H \in \mathcal{H}$ and $\sigma \in \mathcal{S}(H)$ equitable. Furthermore, let Λ_g be the set of all pairs (H, σ) with $H \in \mathcal{H}$ and σ a good 2-coloring of H . Then $\mathbb{E}[Z_e] = |\Lambda_e|/|\mathcal{H}|$ and $\mathbb{E}[Z_g] = |\Lambda_g|/|\mathcal{H}|$. Hence, it suffices to show that $|\Lambda_e| \sim |\Lambda_g|$. But this is evident from Proposition 4.6. Indeed, Proposition 4.6 implies that $|\{H \in \mathcal{H} : (H, \sigma) \in \Lambda_g\}| \sim |\{H \in \mathcal{H} : (H, \sigma) \in \Lambda_e\}|$ for any equitable σ . \square

Corollary 4.8 *Suppose that $r < r_{\text{cond}}$. For any equitable σ we have*

$$\mathbb{P}[\sigma \text{ is a good 2-coloring of } H_k(n, m)] \sim \mathbb{P}[\sigma \text{ is a valid 2-coloring of } H_k(n, m)].$$

Proof. Since the total number of equitable $\tau \in \{0, 1\}^V$ equals $2^{\binom{n}{n/2}}$, and because the uniform distribution over hypergraphs is invariant under permutations of the vertices, we have

$$\begin{aligned} \mathbb{E}[Z_g] &= 2^{\binom{n}{n/2}} \mathbb{P}[\sigma \text{ is a good 2-coloring of } H_k(n, m)], \\ \mathbb{E}[Z_e] &= 2^{\binom{n}{n/2}} \mathbb{P}[\sigma \text{ is a 2-coloring of } H_k(n, m)]. \end{aligned}$$

Hence, the assertion follows from Corollary 4.7. \square

We are going to compute the second moment $\mathbb{E}[Z_g^2]$. The exact same calculation that we used to prove Fact 4.2 shows that $\mathbb{E}[Z_g^2] \leq C \cdot \mathbb{E}[Z_g]^2$ if for any equitable σ we have

$$\mathbb{E}[Z_g | \sigma \text{ is a good 2-coloring}] \leq C \cdot \mathbb{E}[Z_g]. \quad (12)$$

Thus, we are left to verify that for $r < r_{\text{cond}}$ there is $C = C(k, r)$ such that (12) holds.

Let $\alpha = 2^{-k/2}$. Letting $Z_g(d)$ denote the number of good 2-colorings at Hamming distance d from σ , we obtain

$$\mathbb{E} \left[\sum_{0 \leq d \leq \alpha n} Z_g(d) | \sigma \text{ is good} \right] \leq \mathbb{E}[|\mathcal{C}(\sigma)| | \sigma \text{ is good}] \leq \mathbb{E}[Z_e]; \quad (13)$$

the last inequality follows because if σ is good, then $\mathcal{C}(\sigma) \leq \mathbb{E}[Z_e]$ with certainty. Further, by Corollary 4.8,

$$\begin{aligned} \sum_{\alpha n < d \leq n/2} \mathbb{E}[Z_g(d) | \sigma \text{ is good}] &\leq \sum_{\alpha n < d \leq n/2} \mathbb{E}[Z_e(d) | \sigma \text{ is good}] \\ &\leq \sum_{\alpha n < d \leq n/2} \mathbb{E}[Z_e(d) | \sigma \text{ is a valid 2-coloring}] \cdot \frac{\mathbb{P}[\sigma \text{ is a valid 2-coloring}]}{\mathbb{P}[\sigma \text{ is good}]} \\ &\sim \sum_{\alpha n < d \leq n/2} \mathbb{E}_{H_k(n, m, \sigma)}[Z_e(d)] \\ &= \Theta(1/n) \sum_{\alpha n < d \leq n/2} \exp[\psi(d/n)] \quad [\text{by Fact 4.3}]. \end{aligned}$$

Furthermore, Lemma 3.3, Lemma 4.5 and Corollary 4.7 imply that

$$\sum_{\alpha n < d \leq n/2} \mathbb{E}[Z_g(d) | \sigma \text{ is good}] \leq (1 + o(1)) \sum_{\alpha n < d \leq n/2} \mathbb{E}[Z_e(d) | \sigma \text{ is a valid 2-coloring}] \leq C' \cdot \mathbb{E}[Z_e] \quad (14)$$

for a certain constant $C' = C'(k, r)$. Since furthermore $Z_g(d) = Z_g(n - d)$ due to the symmetry of the 2-coloring problem with respect to swapping the color classes, (13) and (14) yield (12) with $C = 2(C' + 1)$.

Hence, we have shown that $\mathbb{E}[Z_g^2] \leq C \cdot \mathbb{E}[Z_g]^2$ for all $r < r_{cond}$. Corollary 4.7 and the Paley-Zygmund inequality therefore imply that

$$\mathbb{P}[Z > 0] \geq \mathbb{P}[Z > \mathbb{E}[Z]/3] \geq \mathbb{P}[Z_g > \mathbb{E}[Z_g]/2] \geq 1/(4C). \quad (15)$$

In particular, the threshold r_{col} for 2-colorability cannot be smaller than r_{cond} , whence indeed $H_k(n, m)$ is 2-colorable w.h.p. for any $r < r_{cond}$. The second claim (3) follows from (15) together with Lemma 3.4.

4.3 Beyond the condensation transition: proof of Theorem 1.3

The goal in this section is establish Theorem 1.3, i.e., to prove that there is a non-empty regime of densities $r_{cond} < r < r_{col}$ in which $H_k(n, m)$ is 2-colorable but $\ln Z < \ln \mathbb{E}[Z] - \Omega(n)$ w.h.p. To get an intuition why this should be the case, consider a density $r > r_{cond} + \varepsilon_k$. In Proposition 4.9 below we will see that w.h.p. (for suitable ε_k), the size $|\mathcal{C}(\sigma)|$ of the ‘local cluster’ in the planted model $H_k(n, m, \sigma)$ is bigger by an exponential factor $\exp(\Omega(n))$ than the expected number $\mathbb{E}[Z]$ of 2-colorings of $H_k(n, m)$. However, if it was true that $\ln Z \sim \ln \mathbb{E}[Z]$, then Proposition 1.5 would imply that the planted model and the Gibbs distribution (first choose $H_k(n, m)$ and then choose $\sigma \in \mathcal{S}(H_k(n, m))$ randomly) are ‘close’. In particular, in a random pair (H, σ) chosen from the Gibbs distribution the local cluster $\mathcal{C}(\sigma)$ should have size $\geq \mathbb{E}[Z] \exp(\Omega(n))$. This would lead to the absurd conclusion that under the Gibbs distribution $Z \geq |\mathcal{C}(\sigma)| \geq \mathbb{E}[Z] \exp(\Omega(n))$ w.h.p. (in obvious contradiction to Markov’s inequality). Hence, intuitively the condensation transition occurs because the size of the local cluster in the planted model $H_k(n, m, \sigma)$ surpasses the expected number $\mathbb{E}[Z]$ of 2-colorings of $H_k(n, m)$. Indeed, it is not difficult to turn this intuition into a proof of part 2 of Theorem 1.3 (see the proof of Theorem 1.3 below). But the above still allows for the possibility that the condensation phase may just be empty, i.e., that the typical size of the local cluster in the planted model $H_k(n, m, \sigma)$ is bounded by $\mathbb{E}[Z]$ for the *entire* regime of r where $H_k(n, m)$ is 2-colorable w.h.p. The purpose of this section is to show that that is not so.

To prove this, we are going to show that w.h.p. $H_k(n, m)$ has a 2-coloring σ whose local cluster $\mathcal{C}(\sigma)$ is smaller than $\mathbb{E}[Z]$, i.e., *much* smaller than the local cluster in the planted model. As we will see in Section 5 below, the size of the local cluster of a 2-coloring σ is governed by the edges that contain precisely one vertex v with color i and $k-1$ vertices with color $1-i$ (with either $i=0$ or $i=1$). Let us call such edges *critical* under σ . Intuitively, critical edges ‘freeze’ v by preventing v from switching to the opposite color $1-i$, thereby reducing the entropy of the local cluster.

Given this intuition, it seems natural to assume that 2-colorings that have a particularly high number of critical edges should have rather small local clusters. Thus, we say that a 2-coloring σ of $H_k(n, m)$ is $(1+\beta)$ -critical if σ is equitable and the total number of critical edges equals $(1+\beta)km/(2^{k-1}-1)$. Let $Z_{1+\beta}$ be the number of $(1+\beta)$ -critical 2-colorings. Furthermore, let us call a $(1+\beta)$ -critical 2-coloring σ *good* if indeed the local cluster

$$\mathcal{C}(\sigma) = \left\{ \tau \in \mathcal{S}(H) : \text{dist}(\sigma, \tau) \leq 2^{-k/2} \right\}$$

satisfies $|\mathcal{C}(\sigma)| \leq \mathbb{E}[Z_{1+\beta}]$, and let $Z_{g,1+\beta}$ be the number of good $(1+\beta)$ -critical 2-colorings.

Proposition 4.9 *For any $k \geq k_0$ there exist a density $r_{crit} > r_{cond}$, $\delta_k > 0$, and $\beta_k > 0$ such that for $r = r_{cond}$ the following three statements hold.*

1. We have $\mathbb{E}[Z_{g,1+\beta_k}] \sim \mathbb{E}[Z_{1+\beta_k}] = \exp(\Omega(n))$.
2. Let $\sigma \in \{0,1\}^V$ be equitable and let $H = H_k(n, m, \sigma)$ be a hypergraph chosen from the planted model. Then w.h.p. the local cluster $\mathcal{C}(\sigma) = \left\{ \tau \in \mathcal{S}(H) : \text{dist}(\sigma, \tau) \leq 2^{-k/2} \right\}$ has size $|\mathcal{C}(\sigma)| > \mathbb{E}[Z] \cdot \exp(\delta_k n)$.

We defer the proof of Proposition 4.9 to Section 5.

In the sequel, we fix $k \geq k_0$ big enough and let $r = r_{crit}$ and $\beta = \beta_k$ be as in Proposition 4.9. In the rest of this section, we are going to carry out a second moment argument for $Z_{g,1+\beta}$ to show the following.

Proposition 4.10 *With r, β as above, we have*

$$\mathbb{E}[Z_{g,1+\beta}^2] \leq C \cdot \mathbb{E}[Z_{g,1+\beta}]^2$$

for some constant $C = C(k) > 0$.

As before, this amounts to showing that

$$\mathbb{E}[Z_{g,1+\beta}|\sigma \text{ is a good } (1+\beta)\text{-critical 2-coloring}] \leq C \cdot \mathbb{E}[Z_{g,1+\beta}], \quad (16)$$

for some number $C = C(k, r) > 0$. To establish (16), we let $Z_{g,1+\beta}(d)$ signify the number of good $(1+\beta)$ -critical 2-colorings at Hamming distance d from σ . Let $\gamma = 2^{-k/2}$. The very definition of ‘good’ ensures that

$$\sum_{d \leq \gamma n} \mathbb{E}[Z_{g,1+\beta}(d)|\sigma \text{ is a good } (1+\beta)\text{-critical 2-coloring}] \leq \mathbb{E}[Z_{g,1+\beta}]. \quad (17)$$

The following bound covers ‘intermediate’ distances.

Lemma 4.11 *We have*

$$\sum_{\gamma n < d < (1/2 - \gamma)n} \mathbb{E}[Z_{g,1+\beta}(d)|\sigma \text{ is a good } (1+\beta)\text{-critical 2-coloring}] = o(1).$$

Proof. Let $\gamma n < d < (1/2 - \gamma)n$. Let τ be equitable and at distance d from σ . Let us briefly say τ is *valid* if τ is a 2-coloring of $H_k(n, m)$. Then

$$\begin{aligned} & \mathbb{P}[\tau \text{ is valid} | \sigma \text{ is a good } (1+\beta)\text{-critical 2-coloring}] \\ &= \frac{\mathbb{P}[\tau \text{ is valid, } \sigma \text{ is a good } (1+\beta)\text{-critical 2-coloring}]}{\mathbb{P}[\sigma \text{ is a good } (1+\beta)\text{-critical 2-coloring}]} \\ &\leq \frac{\mathbb{P}[\sigma, \tau \text{ are valid}]}{\mathbb{P}[\sigma \text{ is a good } (1+\beta)\text{-critical 2-coloring}]} \\ &= \mathbb{P}[\tau \text{ is valid} | \sigma \text{ is valid}] \cdot \frac{\mathbb{P}[\sigma \text{ is valid}]}{\mathbb{P}[\sigma \text{ is a good } (1+\beta)\text{-critical 2-coloring}]} \\ &= \mathbb{P}[\tau \text{ is valid} | \sigma \text{ is valid}] \cdot \frac{\mathbb{E}[Z]}{\mathbb{E}[Z_{g,1+\beta}]} \leq \mathbb{P}[\tau \text{ is valid} | \sigma \text{ is valid}] \cdot \mathbb{E}[Z], \end{aligned}$$

because $\mathbb{E}[Z_{g,1+\beta}] > 1$ by our choice of β . Hence, Lemma 4.5 yields

$$\begin{aligned} \ln \mathbb{E}[Z_{g,1+\beta}(d)|\sigma \text{ is a good } (1+\beta)\text{-critical 2-coloring}] &\leq \ln \mathbb{E}[Z] + \ln \mathbb{E}[Z_e(d)] \\ &\leq -\psi(1/2) + \psi(d/n) + o(1) < 0, \end{aligned}$$

as $\gamma < d/n < 1/2 - \gamma$. Summing over d yields the assertion. \square

Thus, we are left to estimate the contribution of distances $(1/2 - \gamma)n \leq d \leq n/2$. We need to characterize the conditional distribution of $H_k(n, m)$ given that some equitable σ is a $(1+\beta)$ -critical 2-coloring. But since $H_k(n, m)$ is just a uniformly random hypergraph with m edges, this is straightforward: let $H_k(n, m_1, m_2, \sigma)$ denote the random hypergraph generated as follows:

- Choose a set E_1 of m_1 edges that are critical with respect to σ uniformly at random.
- Choose a set E_2 of m_2 edges that are bichromatic under σ but not critical uniformly at random.
- Let $H_k(n, m_1, m_2, \sigma) = (V, E_1 \cup E_2)$.

Then for $m_1 = (1 + \beta)km/(2^{k-1} - 1)$ and for $m_2 = m - m_1$, the conditional distribution of $H_k(n, m)$ given that σ is $(1 + \beta)$ -critical is precisely $H_k(n, m_1, m_2, \sigma)$.

To estimate $\mathbb{E}_{H_k(n, m_1, m_2, \sigma)} [Z_{g, 1+\beta}(d)]$, we need to study the conditional probability that a certain equitable τ at distance d from σ is $(1 + \beta)$ -critical.

Lemma 4.12 *Let τ be equitable at distance $d = \alpha n$ from σ . Then $\mathbb{P}_{H_k(n, m_1, m_2, \sigma)} [\tau \text{ is } 1 + \beta\text{-critical}] \sim \mathcal{E}(\alpha)$, where*

$$\begin{aligned} \mathcal{E}(\alpha) &= (1 - v_1)^{m_1} (1 - v_2)^{m_2} \mathbb{P} [\text{Bin}(m_1, u_1/(1 - v_1)) + \text{Bin}(m_2, u_2/(1 - v_2)) = m_1] \text{ with} \\ u_1 &= (1 - \alpha)^k + \alpha^k + (k - 1)\alpha^2(1 - \alpha)^{k-2} + (k - 1)\alpha^{k-2}(1 - \alpha)^2, \\ v_1 &= \alpha(1 - \alpha)^{k-1} + (1 - \alpha)\alpha^{k-1}, \\ u_2 &= \frac{k(1 - \alpha^k - (1 - \alpha)^k - \alpha^{k-1}(1 - \alpha) - \alpha(1 - \alpha)^{k-1} - (k - 1)\alpha^{k-2}(1 - \alpha)^2 - (k - 1)\alpha^2(1 - \alpha)^{k-2})}{2^{k-1} - k - 1}, \\ v_2 &= \frac{1 - 2[\alpha^k + (1 - \alpha)^k + 2k\alpha(1 - \alpha)^{k-1} + 2k\alpha^{k-1}(1 - \alpha)]}{2^k - 2k - 2}. \end{aligned}$$

Proof. By enumerating all possibilities, we see that the probability that a given edge that is critical under σ also is critical under τ equals u_1 (either all its vertices have the same color under both σ and τ , or they have opposite colors under τ , or the colors of the supporting vertex and exactly one other vertex differ, or the supporting vertex has the same color and the colors of exactly $k - 2$ others differ). Similarly, the probability that an edge that is critical under σ is monochromatic under τ works out to be v_1 .

Now, take a random edge that has $2 \leq l \leq k - 2$ vertices of color 1 under σ . The probability the edge is monochromatic under τ equals

$$\alpha^l(1 - \alpha)^{k-l} + (1 - \alpha)^l\alpha^{k-l}.$$

Convoluting this formula with the distribution of the number of edges with a given number of vertices of color one under σ , we obtain

$$v_2 = \sum_{l=2}^{k-2} \frac{\binom{k}{l}(\alpha^l(1 - \alpha)^{k-l} + (1 - \alpha)^l\alpha^{k-l})}{2^k - 2k - 2}.$$

This is the probability that a random edge that is neither critical nor monochromatic under σ is monochromatic under τ .

Furthermore, the probability that a random edge that has precisely l vertices of color 1 is critical under τ equals

$$l\alpha^{l-1}(1 - \alpha)^{k-l+1} + l(1 - \alpha)^{l-1}\alpha^{k-l+1} + (k - l)(1 - \alpha)^{l+1}\alpha^{k-l-1} + (k - l)\alpha^{l+1}(1 - \alpha)^{k-l-1}.$$

Convoluting this formula with the distribution of the number of edges with a given number of vertices of color one under σ , we get

$$u_2 = \sum_{l=2}^{k-2} \binom{k}{l} \frac{l\alpha^{l-1}(1 - \alpha)^{k-l+1} + l(1 - \alpha)^{l-1}\alpha^{k-l+1} + (k - l)(1 - \alpha)^{l+1}\alpha^{k-l-1} + (k - l)\alpha^{l+1}(1 - \alpha)^{k-l-1}}{2^k - 2k - 2}.$$

This is the probability that a random edge that is neither critical nor monochromatic under σ is critical under τ . The *conditional* probability of a random edge being bichromatic and critical resp. not critical under σ is thus

$$\frac{u_1}{1-v_1} \text{ resp. } \frac{u_2}{1-v_2}.$$

Since the m edges are drawn independently up to the trivial dependence that no edge is drawn twice, we thus see that the probability that τ is $1 + \beta$ -critical is $(1 + o(1))\mathcal{E}(\alpha)$. \square

Corollary 4.13 *For any $0 < d < n$ we have*

$$\begin{aligned} \frac{1}{n} \ln \mathbb{E}_{H_k(n, m_1, m_2, \sigma)} [Z_{g, 1+\beta}(d)] &\sim g(d/n), & \text{with} \\ g(\alpha) &= h(\alpha) + \frac{1}{n} \ln \mathcal{E}(\alpha), \text{ where } h(x) = -x \ln(x) - (1-x) \ln(1-x). \end{aligned}$$

Proof. This simply follows from Lemma 4.12 and the fact that the number of τ at distance d from σ is $\binom{n}{d}$ and $\frac{1}{n} \ln \binom{n}{d} \sim h(d/n)$ by Stirling. \square

Lemma 4.14 *The function g from Corollary 4.13 takes its unique maximum in the interval $(1/2 - \gamma, 1/2 + \gamma)$ at $1/2$, and $g''(1/2) < 0$.*

Proof. Let τ be equitable and at distance αn from σ . Moreover, let $X_1(\alpha)$ be the number of edges of $H_k(n, m_1, m_2, \sigma)$ that are critical under both σ, τ . In addition, let $X_2(\alpha)$ be the number of edges that are critical under τ but not under σ . As $H_k(n, m_1, m_2, \sigma)$ consists of two independent ‘portions’ of random edges, namely m_1 that are critical under σ and another m_2 that are not, $X_1(\alpha), X_2(\alpha)$ are independent. Furthermore,

$$X_1(\alpha) \sim \text{Bin}(m_1, q_1(\alpha)), \quad X_2(\alpha) \sim \text{Bin}(m - m_1, q_2(\alpha)),$$

where $q_1(\alpha) = \frac{u_1}{1-v_1}(\alpha)$, $q_2(\alpha) = \frac{u_2}{1-v_2}(\alpha)$. Let

$$p(\alpha) = \mathbb{P}[X_1(\alpha) + X_2(\alpha) = m_1], \quad p(\alpha, x_1, x_2) = \mathbb{P}[X_1(\alpha) = x_1 \wedge X_2(\alpha) = x_2],$$

so that

$$p(\alpha) = \sum_{x_1, x_2: x_1 + x_2 = m_1} p(\alpha, x_1, x_2).$$

Let us first investigate the point $\alpha = 1/2$. As $q_1(1/2) = q_2(1/2)$, we have

$$(X_1 + X_2)(1/2) \sim \text{Bin}(m_1, q).$$

with

$$q = q_1(1/2) = q_2(1/2) = \frac{2k}{2^k - 2} = \frac{k}{2^{k-1} - 1}.$$

Using Lemma 3.2, we can compute $p(1/2)$ directly:

$$\begin{aligned} p(1/2) &= \binom{m}{m_1} q^{m_1} (1-q)^{m-m_1} \\ &= \Theta(n^{-1/2}) \cdot \exp \left[-mq \cdot \varphi \left(\frac{m_1 - mq}{mq} \right) - (1-q)m \cdot \varphi \left(\frac{mq - m_1}{m(1-q)} \right) \right], \end{aligned}$$

where $\varphi(x) = (1+x) \ln(1+x) - x$. Hence,

$$p(1/2) = \Theta(n^{-1/2}) \cdot \exp \left[-m \left(q \cdot \varphi(\beta) - (1-q) \varphi \left(\frac{\beta q}{q-1} \right) \right) \right].$$

To proceed, we need to decompose this expression according to the individual contributions of $X_1(1/2)$, $X_2(1/2)$. Let x_1, x_2 be such that $x_1 + x_2 = m_1$. Since $X_1(1/2)$, $X_2(1/2)$ are independent, we have

$$\begin{aligned} p(1/2, x_1, x_2) &= \mathbb{P}[X_1 = x_1 \wedge X_2 = x_2] \\ &= \mathbb{P}[\text{Bin}(m_1, q) = x_1] \cdot \mathbb{P}[\text{Bin}(m_2, q) = x_2] \\ &= \Theta(n^{-1}) \exp \left[-qm_1\varphi \left(\frac{x_1 - m_1q}{m_1q} \right) - (1-q)m_1\varphi \left(\frac{m_1q - x_1}{(1-q)m_1} \right) \right] \\ &\quad \cdot \exp \left[-qm_2\varphi \left(\frac{x_2 - m_2q}{m_2q} \right) - (1-q)m_2\varphi \left(\frac{m_2q - x_2}{(1-q)m_2} \right) \right], \end{aligned}$$

and

$$p(1/2) = \sum_{x_1+x_2=m_1} p(1/2, x_1, x_2).$$

Similarly, for general α we have

$$\begin{aligned} p(\alpha, x_1, x_2) &= \Theta(n^{-1}) \exp \left[-q_1(\alpha)m_1\varphi \left(\frac{x_1 - m_1q_1(\alpha)}{m_1q_1(\alpha)} \right) - (1-q_1(\alpha))m_1\varphi \left(\frac{m_1q_1(\alpha) - x_1}{(1-q_1(\alpha))m_1} \right) \right] \\ &\quad \cdot \exp \left[-q_2(\alpha)m_2\varphi \left(\frac{x_2 - m_2q_2(\alpha)}{m_2q_2(\alpha)} \right) - (1-q_2(\alpha))m_2\varphi \left(\frac{m_2q_2(\alpha) - x_2}{(1-q_2(\alpha))m_2} \right) \right], \end{aligned}$$

and

$$p(\alpha) = \sum_{x_1+x_2=m_1} p(\alpha, x_1, x_2).$$

As $u_i(1-\alpha) = u_i(\alpha)$ for $\alpha \in (0, 1)$, we have $u'_i(1/2) = v'_i(1/2) = 0$ for $i = 1, 2$. Moreover, a direct calculation shows that

$$\begin{aligned} u''_1(1/2) &= O(k^2/2^k), \\ u''_2(1/2) &= O(k^3/4^k). \end{aligned}$$

Hence, by the chain rule,

$$\frac{\partial^2}{\partial \alpha^2} (\ln p(\alpha, x_1, x_2)) \big|_{\alpha=1/2} = O(k^3/2^k),$$

whence

$$\ln \frac{p(1/2 - \delta, x_1, x_2)}{p(1/2, x_1, x_2)} = \delta^2 \cdot O(k^3/2^k) + O(\delta^3). \quad (18)$$

Furthermore, as $v''_1(1/2) = O(k^2/2^k)$, $v''_2 = O(k^3/4^k)$, we also obtain

$$\ln \frac{(1 - v_1(1/2 - \delta))^{m_1} (1 - v_2(1/2 - \delta))^{m_2}}{(1 - v_1(1/2 - \delta))^{m_1} (1 - v_2(1/2 - \delta))^{m_2}} = \delta^2 \cdot O(k^3/2^k) + O(\delta^3). \quad (19)$$

As the derivatives of the entropy are

$$\begin{aligned} h'(\alpha) &= -\ln \alpha + \ln(1 - \alpha), \\ h''(\alpha) &= -\frac{1}{\alpha} - \frac{1}{1 - \alpha}, \end{aligned}$$

(18) and (19) yield

$$\mathcal{E}(1/2 - \delta) = -4\delta^2 + \delta^2 \cdot O(k^3/2^k) + O(\delta^3). \quad (20)$$

Finally, (20) shows that g takes its unique maximum in $(1/2 - \gamma, 1/2 + \gamma)$ at $1/2$, and $g''(1/2) < 0$. \square

Combining (17), Lemma 4.11, and Lemma 4.14, we obtain (16). This completes the proof of Proposition 4.10.

Proof of Theorem 1.3. Let $r = r_{crit}$, $\delta = \delta_k$, $\beta = \beta_k$ be as in Proposition 4.9. Then by Proposition 4.10, the probability that $H_k(n, m)$ is 2-colorable is bounded from below by a positive constant. As 2-colorability in $H_k(n, m)$ has a sharp threshold, this implies that $r_{col} > r_{crit}$. This proves the first assertion of Theorem 1.3.

To prove the second assertion, fix $r = r_{crit}$. Then the second part of Proposition 4.9 implies that w.h.p.

$$\frac{1}{n} \ln Z(H_k(n, m, \sigma)) > \frac{1}{n} \ln \mathbb{E}[Z] + \delta n. \quad (21)$$

However, there is no obvious way to derive the second assertion in Theorem 1.3 directly from (21), because it is not clear *a priori* that the random variable $\frac{1}{n} \ln Z(H_k(n, m, \sigma))$ is tightly concentrated. Therefore, we will replace it by another random variable for which concentration is easy to show. Namely, for any $b > 0$ we let

$$\mathcal{Z}_b = \sum_{\sigma \in \{0,1\}^n} \exp(-bw(\sigma)),$$

where $w(\sigma)$ is the number of monochromatic edges under σ . (The above random variable is called the partition function at inverse temperature b .) The random variable $\frac{1}{n} \ln \mathcal{Z}_b$ satisfies a Lipschitz condition: either adding or removing a single edge can change the value of $\frac{1}{n} \ln \mathcal{Z}_b$ by at most b . Therefore, Azuma's inequality implies that in both the random hypergraph $H_k(n, m)$ and in the planted model $H_k(n, m, \sigma)$ we have

$$\mathbb{P}[|\ln \mathcal{Z}_b - \mathbb{E} \ln \mathcal{Z}_b| > y] \leq 2 \exp\left[-\frac{y^2}{2m}\right]. \quad (22)$$

We are going to derive an upper bound on $\mathbb{E} \ln \mathcal{Z}_b(H_k(n, m))$. To this end, let S_μ denote the number of $\sigma \in \{0, 1\}^n$ with $w(\sigma) = \mu$ in $H_k(n, m)$. Then

$$\mathcal{Z}_b = \sum_{\mu=0}^m \exp(-b\mu) \cdot |S_\mu|. \quad (23)$$

Furthermore, letting $\mu = \gamma m$ for some small $\gamma > 0$ and using Lemma 3.2, we obtain

$$\begin{aligned} \frac{1}{n} \ln \mathbb{E} S_\mu &\sim \ln 2 + \frac{1}{n} \ln \mathbb{P}[\text{Bin}(m, 2^{1-k}) = \mu] \\ &= \ln 2 - \frac{r}{2^{k-1}} \left[\varphi(1 - 2^{k-1}\gamma) + (2^{k-1} - 1) \varphi\left(-\frac{2^{k-1}\gamma - 1}{2^{k-1} - 1}\right) \right] \end{aligned} \quad (24)$$

where $\varphi(x) = (1+x) \ln(1+x) - x$. Plugging (24) into (23), we see that

$$\frac{1}{n} \mathbb{E} \ln \mathcal{Z}_b \leq \frac{1}{n} \ln \mathbb{E} Z(H_k(n, m)) + \varepsilon_b, \quad (25)$$

where $\varepsilon_b \rightarrow 0$ as $b \rightarrow \infty$. Intuitively, this mirrors the fact that the partition function is dominated by assignments that violate an $o_b(1)$ -fraction of all clauses as $b \rightarrow \infty$. From now on, fix b large enough so that $\varepsilon_b < \delta/4$. Thus, (22) and (25) imply that

$$\mathbb{P}_{H_k(n, m)} \left[\frac{1}{n} \ln \mathcal{Z}_b > \frac{1}{n} \ln \mathbb{E}[Z(H_k(n, m))] + \frac{\delta}{3} \right] = o(1). \quad (26)$$

We will contrast (26) with the situation in the planted model. Let $\xi > 0$ be sufficiently small and let $\tau \in \{0, 1\}^n$ be such that $||\tau^{-1}(0)| - |\tau^{-1}(1)|| < \xi n$. Then there is an equitable σ such that $\text{dist}(\sigma, \tau) \leq \xi n$.

In $H_k(n, m, \sigma)$ the number of edges that are monochromatic under τ has a binomial distribution with mean $\leq k\xi m$. Therefore, we obtain

$$\frac{1}{n} \mathbb{E} \ln \mathcal{Z}_b(H_k(n, m, \tau)) \geq \frac{1}{n} \mathbb{E} \ln \mathcal{Z}_b(H_k(n, m, \sigma)) - km\xi - o(1).$$

Combining this with (21) and choosing $\xi > 0$ sufficiently small, we thus get

$$\frac{1}{n} \mathbb{E} \ln \mathcal{Z}_b(H_k(n, m, \tau)) \geq \frac{1}{n} \ln \mathbb{E} Z(H_k(n, m)) + 2\delta/3.$$

Hence, (22) yields that

$$\mathbb{P} \left[\frac{1}{n} \ln \mathcal{Z}_b(H_k(n, m, \tau)) < \frac{1}{n} \ln \mathbb{E} Z(H_k(n, m)) + \delta/2 \right] \leq \exp(-\xi'' n), \quad (27)$$

with $\xi'' > 0$.

To complete the proof, consider the set Λ of all pairs (H, τ) of hypergraphs H on $V = \{1, \dots, n\}$ with m edges and 2-colorings τ . Let

$$\begin{aligned} \Lambda' &= \left\{ (H, \tau) \in \Lambda : \frac{1}{n} \ln \mathcal{Z}_b(H) < \frac{1}{n} \ln \mathbb{E} Z(H_k(n, m)) + \delta/2 \right\}, \\ \Lambda'' &= \left\{ (H, \tau) \in \Lambda : \left| |\tau^{-1}(1)| - |\tau^{-1}(0)| \right| \geq \xi n \right\}. \end{aligned}$$

Clearly, (27) shows that

$$|\Lambda' \setminus \Lambda''| \leq \exp(-\xi'' n) |\Lambda|.$$

Furthermore, since the number of hypergraphs H for which $\tau \in \{0, 1\}^n$ is a 2-coloring is maximized for equitable τ , we have $|\Lambda''| \leq \exp(-\xi''' n) |\Lambda|$, with $\xi''' > 0$ sufficiently small. Hence,

$$|\Lambda'| \leq 2 \exp(-\xi''' n) |\Lambda|. \quad (28)$$

Now, suppose that α_1, α_2 are such that

$$\mathbb{P} \left[\frac{1}{n} \ln Z(H_k(n, m)) \geq \frac{1}{n} \ln \mathbb{E} Z(H_k(n, m)) - \alpha_1 n \right] \geq \alpha_2.$$

Since $H_k(n, m)$ is uniformly distributed over all $\binom{n}{m}$ hypergraphs with m edges, we obtain

$$|\Lambda'| \geq \alpha_2 \binom{n}{m} \mathbb{E} [Z(H_k(n, m))] \exp(-\alpha_1 n) \geq \alpha_2 \exp(-\alpha_1 n) \cdot |\Lambda|. \quad (29)$$

Setting $\alpha_1 = \xi'''/2$ and comparing (28) with (29), we see that necessarily $\alpha_2 = o(1)$. This proves (4) in the case that $r = r_{crit}$.

Finally, consider any density $r_{crit} < r < r_{col}$. We generate the random hypergraph $H_k(n, m)$ in two ‘portions’ H_1 and H_2 . Namely, letting $m_1 = r_{cond}n$ and $m_2 = (r - r_{cond})n$, we let $H_1 = H_k(n, m_1)$. Then H_2 is simply obtained by adding another m_2 random edges to H_1 . By the above, we know that w.h.p.

$$Z(H_1) \leq \mathbb{E} [Z(H_1)] \cdot \exp(-\Omega(n)).$$

Furthermore, a new random edge is bichromatic under a 2-coloring of H_1 with probability $1 - 2^{1-k}$, we have

$$\mathbb{E} [Z(H_2)|H_1] \leq Z(H_1) \cdot (1 - 2^{1-k})^{m_2}.$$

Thus, w.h.p.

$$Z(H_k(n, m)) = Z(H_2) \leq \mathbb{E} [Z(H_1)] \cdot \exp(-\Omega(n)) (1 - 2^{1-k})^{m_2} = \mathbb{E} [Z(H_k(n, m))] \exp(-\Omega(n)),$$

as claimed. \square

4.4 Proof of Lemma 4.5

Since $\psi(1-x) = \psi(x)$, we only need to work with $x \leq 1/2$. Let $r = (2^{k-1} - 1)(c/2^k + \ln 2)$ for $|c| \leq 4$. Let $h(x) = -x \ln x - (1-x) \ln(1-x)$. Then

$$\psi(x) \leq h(x) - \frac{r}{2^{k-1} - 1} (1 - x^k - (1-x)^k) \leq h(x) - (c/2^k + \ln 2)(1 - x^k - (1-x)^k).$$

Suppose that $x > 2^{-k/2}$ but $x < 1/(1.01k)$. Then

$$\begin{aligned} \psi(x) &\leq x(1 - \ln x) - (c/2^k + \ln 2) \left(1 - (1-x)^k\right) + 2^{-k} \\ &\leq x(1 - \ln x) - (c/2^k + \ln 2) (kx - (kx)^2) + 2^{-k} \\ &\leq x [1 - \ln x - k(1 - kx) \ln 2] + 2^{3-k} < -2^{3-k} \leq -\psi(1/2), \end{aligned} \quad (30)$$

provided that $k \geq k_0$ is large enough. Furthermore, for $1/(1.01k) < x < 0.49$ we have

$$\begin{aligned} \psi(x) &\leq h(x) - (c/2^k + \ln 2) \left(1 - (1-x)^k\right) + 2^{-k} \\ &\leq h(x) - (c/2^k + \ln 2)(\exp(-kx) - 1) + 2^{-k} \\ &\leq h(x) - (\exp(-kx) - 1) \ln 2 + 2^{3-k} < -2^{3-k} \leq -\psi(1/2), \end{aligned} \quad (31)$$

again for $k \geq k_0$ large enough.

Finally, around $x = 1/2$ we can expand ψ as follows. Since $\psi(1-x) = \psi(x)$, it is clear that $\psi'(1/2) = 0$. Furthermore, $\psi''(1/2) = -4 + o_k(1)$, and $\psi'''(1/2) \leq h'''(1/2) + o_k(1) = o_k(1)$. Therefore, for $k \geq k_0$ large enough we can expand ψ around $1/2$ as

$$\psi\left(\frac{1}{2} - \delta\right) = \psi(1/2) - (4 + o_k(1))\delta^2 + O(\delta^3). \quad (32)$$

Thus, the lemma follows from (30)–(32).

5 The local cluster: proof of Propositions 4.6 and 4.9

5.1 Outline

In this section we prove Propositions 4.6 and 4.9. Fix an equitable 2-coloring $\sigma : V \rightarrow \{0, 1\}$ and recall that an edge e of a hypergraph H is *critical* under σ if there is a color $i \in \{0, 1\}$ and a vertex $v \in E$ such that $\sigma(v) = i$ and $\sigma(w) = 1 - i$ for all $w \in e \setminus \{v\}$. In this case, we say that v *supports* the edge e (under σ).

We are going to study the size of the local cluster in the $H_k(n, m_1, m_2, \sigma)$ model from Section 4.3:

- Choose a set E_1 of m_1 edges that are critical with respect to σ uniformly at random.
- Choose a set E_2 of m_2 edges that are bichromatic under σ but not critical uniformly at random.
- Let $H_k(n, m_1, m_2, \sigma) = (V, E_1 \cup E_2)$.

We are going to expose the edges of $H_k(n, m_1, m_2, \sigma)$ in two portions: let H_1 contain the m_1 critical edges, and let H_2 contain the rest. Let $\lambda = m_1/n$ be the expected number of edges that any one vertex supports.

We will need the following simple expansion property of the random hypergraph H_1 .

Lemma 5.1 *Let $\zeta < 1/3$. W.h.p. H_1 has the following property. Suppose that $S \subset V$ has size $|S| = \zeta n$. Then w.h.p. the total number of edges supported by vertices in S is bounded by $\zeta(e^3\lambda - \ln \zeta)n$.*

Proof. We use a first moment argument. Let $\xi = e^3\lambda - \ln \zeta$ and $\mu = \xi\zeta$. The probability that there is a set S of size ζn that supports a total of μn edges is bounded by

$$\begin{aligned} \binom{n}{\zeta n} \binom{m_1}{\mu n} \zeta^{\mu n} &\leq \left[\left(\frac{e}{\zeta} \right)^\zeta \left(\frac{e\lambda\zeta}{\mu} \right)^\mu \right]^n = \left[\frac{e}{\zeta} \left(\frac{e\lambda\zeta}{\mu} \right)^\xi \right]^{\zeta n} \\ &\leq \left[\frac{e}{\zeta} \left(\frac{e\lambda}{\xi} \right)^\xi \right]^{\zeta n} \leq [e\zeta]^{\zeta n} = o(1), \end{aligned}$$

by our choice of ξ and because $\zeta < 1/3$. □

Lemma 5.2 *Let $l \geq 0$ be fixed. W.h.p. the number of vertices that support precisely l edges is*

$$(1 + o(1))n \cdot \frac{\lambda^l}{l! \exp(\lambda)}$$

Proof. The number of edges that any one vertex supports is binomial with mean λ . Hence, the Poisson approximation to the binomial distribution shows that the probability that some vertex v supports precisely l edges is $(1 + o(1)) \frac{\lambda^l}{l! \exp(\lambda)}$. In effect, letting X_l be the number of vertices with this property, we see that

$$\mathbb{E}X_l = (1 + o(1))n \cdot \frac{\lambda^l}{l! \exp(\lambda)}.$$

Furthermore, X_l satisfies a Lipschitz condition: adding or removing a single edge can change the value of X_l by at most one. Therefore, Azuma's inequality shows that $X_l = (1 + o(1))n \cdot \frac{\lambda^l}{l! \exp(\lambda)}$ w.h.p. □

In particular, Lemma 5.2 shows that the total number of vertices that do not support any edges is $(1 + o(1)) \exp(-\lambda)n$ w.h.p. Now, consider the following construction of a set $U \subset V$.

1. Initially, let U consist of all vertices that do not support any edges.
2. While there is a vertex $v \notin U$ that does not support an edge that does not contain a vertex from U , add v to U .

The above is an adaptation of the ‘whitening process’ from [6] to random hypergraph 2-coloring.

Let H_U be the hypergraph with vertex set U and edge set

$$\{e \cap U : e \in E(H_1), |e \cap U| \geq 2\}.$$

In general, this is going to be a non-uniform hypergraph.

Proposition 5.3 *W.h.p. the set U has size*

$$|U| = n \left[\exp(-\lambda) + \lambda(k-1) \exp(-2\lambda) + O(7.1^{-k}) \right]$$

and enjoys the following properties.

U1. *The set $S_0 \subset U$ of variables that do not support a clause has size $|S_0| = (1 + o(1))n \exp(-\lambda)$.*

U2. *There is a set S_1 of size*

$$(1 + o(1))n \left[\lambda(k-1) \exp(-2\lambda) + O(7^{-k}) \right]$$

such that all vertices in S_1 support exactly one edge that contains precisely one other vertex from U , which indeed belongs to S_0 .

U3. Apart from the edges resulting from **U2**, H_U contains no more than $nO(7.1^{-k})$ further edges.

We defer the proof of Proposition 5.3 to Section 5.2.

We say that $R \subset V$ is *rigid* if for any 2-coloring τ of H_1 such that $\tau(v) \neq \sigma(v)$ we have

$$|\{v \in R : \tau(v) \neq \sigma(v)\}| \geq n/k^3.$$

In Section 5.3 we will prove the following.

Proposition 5.4 *W.h.p. there is a rigid set $R \subset V \setminus U$ of size $|R| \sim |V \setminus U|$.*

We now have sufficient information about the random hypergraph $H_k(n, m_1, m_2, \sigma) = H_1 \cup H_2$ to prove Propositions 4.6 and 4.9.

Proof of Proposition 4.6. Proposition 4.6 deals with the random hypergraph $H_k(n, m, \sigma)$, in which the number of critical edges has a binomial distribution $\text{Bin}(m, k/(2^{k-1} - 1))$. Hence, by Chernoff bounds the number of critical edges is $(1 + o(1))mkr/(2^{k-1} - 1) = (1 + o(1))\lambda n$ w.h.p., with $\lambda = kr/(2^{k-1} - 1)$. Thus, to study $H_k(n, m, \sigma)$ it suffices to investigate $H_k(n, m_1, m_2, \sigma)$ with $m_1 \sim \lambda n$ and $m_2 \sim m - \lambda n$.

To prove Proposition 4.6 we merely need to derive an *upper* bound on the size $|\mathcal{C}(\sigma)|$ of the local cluster. Thus, it suffices to bound the size of the local cluster

$$\mathcal{C}^1(\sigma) = \left\{ \tau : \text{dist}(\sigma, \tau) \leq 2^{-k/2}n, \quad \tau \text{ is a 2-coloring of } H_1 \right\}$$

of H_1 . By Proposition 5.4, we have w.h.p.

$$\frac{1}{n} \ln |\mathcal{C}(\sigma)| \leq \frac{1}{n} \ln |\mathcal{C}^1(\sigma)| = \frac{1}{n} \ln Z(H_U).$$

Hence, we just need to bound the number $Z(H_U)$ of 2-colorings of H_U . By Proposition 5.3 we may assume that H_U has the properties **U1–U3**.

If this is indeed the case, we can estimate $\ln Z_U$ as follows. Let H'_U be the hypergraph obtained from H_U by omitting all edges that are incident with a vertex from $U \setminus (S_0 \cup S_1)$. Then each edge of H'_U has size 2 and contains precisely one vertex from S_1 and one vertex from S_0 . Moreover, each vertex from S_1 is incident with exactly one such edge, and indeed supports this edge under σ . Hence, H'_U is just a collection of stars in which all non-isolated vertices in S_1 are leaves, and therefore the total number of 2-colorings of H'_U is simply equal to $2^{|S_0|}$. Thus,

$$\frac{1}{n} \ln |\mathcal{C}(\sigma)| \leq \frac{1}{n} \ln Z(H_U) \leq \frac{1}{n} \ln Z(H'_U) \leq \frac{|S_0|}{n} \ln 2 \leq (1 + o(1)) \exp(-\lambda) \ln 2.$$

A straightforward computation shows that this is indeed less than $\frac{1}{n} \ln \mathbb{E}[Z_e(H_k(n, m))]$ if $r < (2^{k-1} - 1) \ln 2$. \square

Proof of Proposition 4.9. We start by obtaining an upper bound on the size of $\mathcal{C}(\sigma)$. Let $\lambda = (1 + \beta)kr/(2^{k-1} - 1)$ for some $\beta \leq 1/k$. We first study the size of the local cluster $\mathcal{C}^1(\sigma)$ in H_1 . By the same argument as in the proof of Proposition 4.6 above, w.h.p. we have

$$\frac{1}{n} \ln |\mathcal{C}^1(\sigma)| \leq \frac{1}{n} \ln Z(H'_U) \leq (1 + o(1)) \exp(-\lambda) \ln 2,$$

where H'_U is a collection of stars as above. While clearly $\frac{1}{n} \ln |\mathcal{C}(\sigma)| \leq \frac{1}{n} \ln |\mathcal{C}^1(\sigma)|$, we need a slightly tighter estimate of $|\mathcal{C}(\sigma)|$.

To obtain this estimate, we need to take the edges of H_2 into consideration. Let E'_2 consist of all edges $e \in H_2$ that contain precisely two vertices from $S_0 \setminus N(S_1)$ and in which all vertices in $V \setminus U$ have the

same color under σ . Since H_2 is independent of H_1 , the number of these edges is binomially distributed with mean

$$\frac{|S_0 \setminus N(S_1)|^2}{n^2} \cdot \frac{\binom{k}{2}}{2^{k-1} - 1} \cdot n \cdot m_2 \geq n \binom{k}{2} \exp(-2\lambda) \ln 2 = \mu_2.$$

By Chernoff bounds, we indeed have $|E'_2| \geq (1 - o(1))\mu_2$ w.h.p. Furthermore, the expected number of vertices in S_0 that are incident with two edges from E'_2 is $\leq O(k^4 \exp(-3\lambda))$; as this number satisfies a Lipschitz condition, it is concentrated by Azuma's inequality. Hence, w.h.p. E'_2 contains a subset E''_2 of size

$$|E''_2|/n \geq \binom{k}{2} \exp(-2\lambda) \ln 2 - O(k^4/8^k)$$

such that E''_2 induces a matching in S_0 . By construction, this matching is disjoint from H'_U . Hence, w.h.p.

$$\frac{1}{n} \ln |\mathcal{C}(\sigma)| \leq \frac{1}{n} \ln |\mathcal{C}^1(\sigma)| - |E''_2| \ln 2 \leq \exp(-\lambda) \left[1 - \binom{k}{2} \exp(-\lambda) \ln 2 \right] \ln 2 + O(k^4/8^k). \quad (33)$$

To derive a matching lower bound, notice that Proposition 5.3 implies that all but $O(7.1^{-k})n$ edges of H_1 belong to the matching H'_U w.h.p. Let F_1 be the set of all vertices that are reachable from the edges in $H_1 \setminus H'_U$. Then $|F_1| \leq 4|H_1 \setminus H'_U| \leq O(7.1^{-k})n$ w.h.p. While we cannot say much about the entropy of the vertices in F_1 , it is clear that $H_U - F_1$ is just a matching from $S_1 \setminus F$ to $S_0 \setminus F$. Therefore, w.h.p.

$$\frac{1}{n} \ln |\mathcal{C}^1(\sigma)| \geq |S_0 \setminus F| \ln 2 \geq (\exp(-\lambda) - O(7.1^{-k})) \ln 2.$$

Let E'_3 be the set of all edges $e \in H_2$ that contain at least three vertices from U such that all vertices in $e \setminus U$ have the same color under σ . Then $\mathbb{E}|E'_3| \leq O(k^3/2^k) \cdot (|U|/n)^3 m \leq O(k^3/8^k)n$. Let F_3 be the set of all vertices in H_U that are reachable from $\{v \in U : \exists e \in E'_3 : v \in e\}$. Since $|E'_3|$ is binomially distributed, we have $|F_3| \leq O(k^4/8^k)n$ w.h.p. Furthermore, let E'_2 be as above. Let F'_2 be the set of all vertices in $N(S_1) \cup U \setminus (S_0 \cup F_1 \cup F_3)$ that are incident with an edge of E'_2 . Then $|F'_2|$ is binomially distributed with mean $\leq O(k^2/2^k) \exp(-\lambda) |U \setminus S_0| m \leq O(k^3/8^k)n$ (by Proposition 5.3), and thus w.h.p. $|F'_2| \leq O(k^3/8^k)n$ by Chernoff bounds. In addition, let F''_2 be the set of all vertices in S_0 that are incident with at least two edges from E'_2 . As we saw above, $F''_2 \leq O(k^4/8^k)n$ w.h.p. Let F_2 be the set of all vertices in H'_U that are reachable from $F'_2 \cup F''_2$; since H'_U is a matching, we have $|F_2| \leq 2|F'_2 \cup F''_2| = O(k^4/8^k)n$ w.h.p. Finally, let E''_2 be the set of all edges in E'_2 that do not contain a vertex from F_2 . Then $|E''_2|/n \leq \binom{k}{2} \exp(-2\lambda) \ln 2 - O(k^4/8^k)$ w.h.p.

Now, E''_2 and H'_U simply induce a matching on $(S_0 \cup S_1) \setminus (F_1 \cup F_2 \cup F_3)$, and this matching is disconnected from all other edges of $H_1 \cup H_2$ that are not already 2-colored given the colors assigned to the vertices in $V \setminus U$. Hence, w.h.p. the number of 2-colorings is at least

$$\begin{aligned} \frac{1}{n} \ln |\mathcal{C}(\sigma)| &\geq \frac{1}{n} \ln |\mathcal{C}^1(\sigma)| - |E''_2| \ln 2 - |F_1 \cup F_2 \cup F_3| \ln 2 \\ &\geq \exp(-\lambda) \left[1 - \binom{k}{2} \exp(-\lambda) \ln 2 \right] \ln 2 - O(7.1^{-k}). \end{aligned} \quad (34)$$

To prove the first claim, we need to combine (33) and (34) with a lower bound on the expected number of $(1 + \beta)$ -critical 2-colorings. Let $q = k/(2^{k-1} - 1)$. The probability $\eta_{1+\beta}$ that an equitable σ is a $(1 + \beta)$ -critical 2-coloring of $H_k(n, m)$ satisfies

$$\ln \eta_{1+\beta} \sim m \ln(1 - 2^{1-k}) + \ln \mathbb{P}[\text{Bin}(m, q) = (1 + \beta)qm].$$

Indeed, the first summand accounts for the probability that σ is a 2-coloring, and the second summand is the probability that given that σ is a 2-coloring, the number of critical edges equals $(1 + \beta)qm$. By the Lemma 3.2, for sufficiently small $\beta > 0$ we have

$$\frac{1}{n} \ln \mathbb{P} [\text{Bin}(m, q) = (1 + \beta)qm] \geq -\frac{\beta^2 qm}{n} \geq -k\beta^2.$$

Hence,

$$\frac{1}{n} \ln \mathbb{E} [Z_{1+\beta}] \geq \frac{1}{n} \ln \mathbb{E} Z - k\beta^2.$$

If $r = 2^{k-1} \ln 2 - c$, then a direct computation shows that

$$\frac{1}{n} \ln \mathbb{E} [Z] = \ln 2 + r \ln (1 - 2^{1-k}) \geq \frac{(2c - \ln 2)}{2^k} - O(4^{-k}). \quad (35)$$

Consequently,

$$\frac{1}{n} \ln \mathbb{E} [Z_{1+\beta}] \geq \frac{(2c - \ln 2)}{2^k} - O(4^{-k}) - k\beta^2. \quad (36)$$

Choose c (and thus r) such that with $\lambda_0 = kr/(2^{k-1} - 1)$ we have

$$\Xi = \exp(-\lambda_0) \left[1 - \binom{k}{2} \exp(-\lambda_0) \ln 2 \right] \ln 2 - 7^{-k} = \frac{1}{n} \ln \mathbb{E} Z + 16^{-k}. \quad (37)$$

A straight computation using (35) shows that $c = \ln 2 + o_k(1)$. Furthermore, (34) and (37) show that for this r w.h.p. in the planted model $H_k(n, m, \sigma)$ the local cluster $\mathcal{C}(\sigma)$ has size $|\mathcal{C}(\sigma)| > \exp(\Omega(n))\mathbb{E}Z$. Let

$$f(\beta) = \exp(-(1 + \beta)\lambda_0) \left[1 - \binom{k}{2} \exp(-(1 + \beta)\lambda_0) \ln 2 \right] \ln 2.$$

Expanding $f(\cdot)$ around $\beta = 0$, we find that

$$f(\beta) - f(0) = -\beta(\exp(-\lambda_0) \ln 2 + O(k^2 4^{-k})) + O(\beta^2)/2^k.$$

Hence, (36) implies that for $\beta^* = 3^{-k}$ we get

$$f(\beta^*) + 7^{-k} < \frac{1}{n} \ln \mathbb{E} [Z_{1+\beta^*}].$$

Further, (33) implies that with $m_1 = (1 + \beta^*)\lambda_0 n$, $m_2 = m - m_1$ in $H_k(n, m_1, m_2, \sigma)$ w.h.p. the local cluster size satisfies $\frac{1}{n} \ln |\mathcal{C}(\sigma)| < \frac{1}{n} \ln \mathbb{E} [Z_{1+\beta^*}]$. This means that r, β^* as above satisfy the conditions in Proposition 4.9. \square

5.2 Proof of Proposition 5.3

Let U_0 be the set of all vertices that do not support any edge. Then w.h.p. $|U_0| \sim n \exp(-\lambda)$ by Lemma 5.2. For each vertex v let $s(v)$ be the number of edges that v supports. Let U_1 be the set of all vertices v with $s(v) \geq 1$ such that all edges supported by v contain a vertex from U_0 .

Lemma 5.5 *W.h.p. we have the following.*

1. The number of vertices v with $s(v) = 1$ such that the edge e supported by v contains exactly one vertex from U_0 is

$$n \left[\lambda(k-1) \exp(-2\lambda) + O(7 \cdot 3^{-k}) \right].$$

2. The number of vertices v with $s(v) = 1$ such that the edge e supported by v contains more than one vertex from $(U_0 \cup U_1) \setminus \{v\}$ is $n \cdot O(7.3^{-k})$.
3. The number of vertices v with $s(v) > 1$ such that all edges e supported by v contain a vertex from U_0 is bounded by $n \cdot O(7.3^{-k})$.

Proof. Let X be the number of vertices as in 1. By Lemma 5.2, the number of vertices v with $s(v) = 1$ is $(1 + o(1))\lambda \exp(-\lambda)n$ w.h.p. Furthermore, given that v satisfies $s(v) = 1$, the $k - 1$ other vertices in the unique edge e that v supports are uniformly distributed over the opposite color class. Hence, again by Lemma 5.2, the number of non-supporting vertices amongst these $k - 1$ vertices has a binomial distribution $\text{Bin}(k - 1, (1 + o(1)) \exp(-\lambda))$ w.h.p. In this case, the probability that exactly one of the $k - 1$ other vertices is non-supporting is $(k - 1) \exp(-\lambda) + O(\exp(-2\lambda))$. Hence, we see that

$$EX = (1 + o(1))\lambda \exp(-\lambda) \cdot [(k - 1) \exp(-\lambda) + O(\exp(-2\lambda))] = n \left[\lambda(k - 1) \exp(-2\lambda) + O(7.9^{-k}) \right].$$

Furthermore, X satisfies $X = EX + o(n)$ w.h.p.; for the number of vertices v with $s(v) = 1$ is concentrated by Lemma 5.2. In addition, for all such v with $\sigma(v) = 0$ the events that the edge e_v supported by v contains a non-supporting vertex are mutually independent. Hence, this number has a binomial distribution and is therefore concentrated by Chernoff bounds (Lemma 3.1). As the same is true of the vertices v with $\sigma(v) = 1$, X is concentrated about its expectation. The other two claims follow from a similar argument. \square

Then the above lemma shows that $|U_1|/n \leq \lambda(k - 1) \exp(-2\lambda) + O(7^{-k})$ w.h.p. Furthermore, the hypergraph $H_{U_1 \cup U_2}$ mostly consists of isolated vertices and edges of size 2 (and only very larger edges).

We now need to analyze how the process for the construction of the set U proceeds. All vertices in $V \setminus (U_0 \cup U_1)$ support at least one edge that does not contain a vertex from U_0 . We will now construct sets $U_j, j \geq 2$, inductively as follows:

let U_j be the set of all vertices $v \in V \setminus \bigcup_{i < j} U_i$ such that all edges supported by v contain a vertex from $\bigcup_{i < j-1} U_i$.

Let $U^* = \bigcup_{j \geq 2} U_j$.

Lemma 5.6 *W.h.p. H_1 has the following property. Let T be a set of size $\leq n/2^{k-2}$. Then the number \hat{T} of critical edges that are supported by a vertex $v \notin T$ but that contain a vertex from T is bounded by $36k^3 2^{-k}n$.*

Proof. We use a first moment argument. Let $t = 2^{2-k}$ and $\mu = 36k^3/2^k$. Then probability of the event described above is bounded by

$$\binom{n}{tn} \binom{m_1}{\mu n} (kt)^{\mu n} \leq \left[\frac{e}{t} \left(\frac{e\lambda kt}{\mu} \right)^{\mu/t} \right]^{tn} \leq \left[2^k \left(\frac{e}{9k} \right)^{9k^3} \right]^{tn} = o(1),$$

as claimed. \square

Lemma 5.7 *W.h.p. we have $|U^*| \leq n \cdot O(7.2^{-k})$.*

Proof. This is based on a branching process argument. More precisely, we consider the following stochastic process. At each time, a vertex can be either alive, neutral, or dead. Initially, all vertices in U_0 are dead, all vertices in U_1 are alive, and all other vertices are neutral. In each round of the process an alive vertex a is chosen arbitrarily (once there is no alive vertex left, the process stops). Every neutral vertex v such that all edges e with $v \in e$ contain either a or a dead vertex is declared alive, and then a is declared dead.

Let A_t be the set of alive vertices after t steps of the process (in particular, $A_0 = U_1$). Let $T_* = |A_0|$, $T^* = 2n \cdot 7.2^{-k}$, and let T be the actual stopping time of the process. The goal is to show that w.h.p.

$$T \leq T_* + T^*,$$

which implies that $U^* \setminus U_1 \leq T^*$.

To prove this bound, we proceed as follows. Consider a time $t \leq T_* + T^*$. There are several ways in which a neutral vertex v can become alive.

Case 1: $s(v) = 1$. By Lemma 5.2 the total number of such vertices is bounded by $(1 + o(1))\lambda \exp(-\lambda)n$ w.h.p. Moreover, v can become alive only if the unique clause that v supports contains a . The probability of this event is bounded by $2k/n$. Hence, the expected number of new alive vertices that arise in this way is $\leq (1 + o(1))2k\lambda \exp(-\lambda)$.

Case 2: $s(v) > 1$ and v has a dead neighbor. By Lemma 5.6 and our assumption on t , the total number of vertices with a dead neighbor is bounded by $36k^3 2^{-k}n$. If v is declared alive at time t , then all edges that contain v but no dead vertex must contain v , and there is at least one such edge. The probability of this event is bounded by $2k/n$. Hence, the expected number of vertices that become alive in this way is $\leq (1 + o(1))72k^4 2^{-k}$.

Case 3: $s(v) > 1$ and v does not have a dead neighbor. In this case all $s(v) \geq 2$ edges that v supports contain a . The probability of this event is $O(n^{-2})$. Hence, the expected number of vertices that become alive in this way is $o(1)$.

Thus, conditioning on the previous history \mathcal{F}_{t-1} of the process, we obtain

$$\mathbb{E}[A_t - A_{t-1} | \mathcal{F}_{t-1}] \leq k^5/2^k.$$

Furthermore, for all neutral v the events that v is activated at time t given \mathcal{F}_{t-1} are mutually independent. Hence, $A_t - A_{t-1}$ given \mathcal{F}_{t-1} is stochastically dominated by a binomial variable B_t with mean $k^5/2^k$. Now, if $T \geq T_* + T^*$, then at least T^* vertices got activated by time $T_* + T^*$, i.e., $\sum_{t=1}^{T_*+T^*} B_t \geq T^*$. Since

$$\mathbb{E} \sum_{t=1}^{T_*+T^*} B_t \leq (T_* + T^*)k^5/2^k < T^*/2,$$

the Chernoff bound from Lemma 3.1 shows that $\mathbb{P} \left[\sum_{t=1}^{T_*+T^*} B_t \geq T^* \right] \leq \exp(-\Omega(n))$. \square

Proof of Proposition 5.3. The above discussion allows us to get a close understanding of the combinatorial structure of the hypergraph H_U . By Lemma 5.7 we have $|U^*| \leq n \cdot O(7.2^{-k})$. Let E^* be the set of all edges supported by a vertex in U^* that contain a vertex in $U_0 \cup U_1$. Lemma 5.1 implies that w.h.p. $|E^*| \leq O(k)|U^*| \leq n \cdot O(7.19^{-k})$. Hence, the set U'_* of all vertices $v \in U_0 \cup U_1$ that occur in an edge from E^* has size $|U'_*| \leq n \cdot O(7.18^{-k})$ w.h.p. Furthermore, let U_* be the set of all vertices $v \in U_0 \cup U_1$ such that either $v \in U_*$ or there is an edge e supported by a vertex in U_1 that contains v and another vertex from $U_0 \cup U_1$, or such that $v \in U_0$ occurs in an edge supported by a vertex $w \in U'_* \cap U_1$. Then by Lemma 5.5 we have $|U_*| \leq nO(7.2^{-k})$.

In summary, we have shown that H_U has the following structure w.h.p.

- The set U_0 of non-supporting variables has size $(1 + o(1))n \exp(-\lambda)$.
- There is a set $U_1 \setminus U_*$ of size

$$(1 + o(1))n \left[\lambda(k-1) \exp(-2\lambda) + O(7.17^{-k}) \right]$$

such that in H_U all vertices in $U_1 \setminus U_*$ support exactly one edge that contains precisely one other vertex from U , which indeed belongs to U_0 .

- Apart from these, H_U contains no more than $nO(7.17^{-k})$ further edges.

This completes the proof of Proposition 5.3. \square

5.3 Proof of Proposition 5.4

As a first step, we will identify a large set of rigid vertices. To this end, we need to say something about the number of edges that the vertices in $V \setminus U$ support. Let $l = 10$.

Lemma 5.8 *W.h.p. the number of vertices $v \notin U$ that support fewer than l edges that do not contain a vertex from U is bounded by $\frac{2\lambda^l}{l!} \exp(-\lambda)n$.*

Proof. By Lemma 5.2 the total number of vertices that support fewer than l edges is $\leq \frac{1.01\lambda^l}{l!} \exp(-\lambda)n$ w.h.p. Moreover, applying Lemma 5.6 to the set U , we see that no more than $36k^3n/2^k < 0.9\frac{\lambda^l}{l!} \exp(-\lambda)n$ critical edges supported by a vertex in $V \setminus U$ contain a vertex from U w.h.p. Each of these edges can create at most one additional vertex in $V \setminus U$ that supports fewer than l edges without a vertex from U . \square

For each $v \notin U$ let $s'(v)$ be the number of edges supported by v that do not contain a vertex from U . By the construction of U , we have $s'(v) \geq 1$ for all $v \notin U$. Furthermore, given the sequence $(s'(v))_{v \in V \setminus U}$, the distribution of the sub-hypergraph of H_1 induced on $V \setminus U$ is very simple: it is obtained by choosing, for each vertex $v \in V \setminus U$ independently, $s'(v)$ edges supported by v and containing a random set of $k-1$ vertices from $V \setminus U$ of color $1 - \sigma(v)$. This follows because the construction of the set U merely imposes the condition that none of the $s'(v)$ remaining edges supported by v contains a vertex from U .

We now decompose the random edges of the sub-hypergraph $H_1 - U$ into two portions. The first portion \mathcal{M} contains for each vertex v one random edge supported by v and containing $k-1$ vertices of color $1 - \sigma(v)$ (with no vertex from U , of course). The second portion \mathcal{H} contains the remaining $s'(v) - 1 \geq 0$ random edges supported by v and containing $k-1$ vertices of color $1 - \sigma(v)$ (again, none of them from U). This decomposition will allow us to construct the desired set R in two independent steps.

The first step is in to find a ‘core’ in the hypergraph \mathcal{H} .

CR1. Initially, let S contain all $v \in V$ that support at least $l/2$ edges.

CR2. While there is $v \in S$ that supports $< l/2$ edges consisting of vertices of S only, remove v from S .

Let $\mathcal{C} = S$ be the final outcome of this process. In order to study $|\mathcal{C}|$, we need the following expansion property of the random hypergraph H_1 .

Lemma 5.9 *W.h.p. the random hypergraph H_1 has the following property. Let $T \subset V$ be a set of size tn with $t \leq 1/(e^2k\lambda)$. Then there are no more than $2tn$ edges that are supported by a vertex in T and that contain a second vertex from T .*

Proof. We use a first moment argument. The probability that there is a set T that violates the above property is bounded by

$$\binom{n}{tn} \binom{m_1}{2tn} (kt^2)^{2tn} \leq \left[\frac{e}{t} \left(\frac{em_1kt^2}{2tn} \right)^2 \right]^{tn} = \left[\frac{e}{t} \left(\frac{e\lambda kt}{2} \right)^2 \right]^{tn} \leq \left(\frac{t}{e} \right)^{tn} = o(1),$$

as claimed. \square

Lemma 5.10 *W.h.p. we have $|V \setminus \mathcal{C}| \leq \lambda^l \exp(-\lambda)n$.*

Proof. Assume that $|V \setminus \mathcal{C}| > \lambda^l \exp(-\lambda)n$. By Lemma 5.8 we may assume that the initial set S contains at least $n(1 - 2\lambda^l \exp(-\lambda)/l!)$ vertices. Hence, if $|V \setminus \mathcal{C}| > \lambda^l \exp(-\lambda)n$, then at some point the process **CR1–CR2** must have removed a set T of size $\lambda^l \exp(-\lambda)n/2$ from the original set S . This set T has the property that each vertex in T supports $l/2 > 2$ edges, each of which must contain another vertex from T . But by Lemma 5.9 no such set T exists w.h.p. \square

Having constructed the set \mathcal{C} , we are now going to ‘attach’ more vertices from $V \setminus U$ to it via the following process.

A1. Let $\mathcal{A}_0 = \mathcal{C}$.

A2. For $t \geq 1$, let \mathcal{A}_t be the set of all vertices $v \in V \setminus U$ such that either $v \in \mathcal{A}_{t-1}$ or the edge $e \in \mathcal{M}$ supported by v has its other $k-1$ vertices in \mathcal{A}_{t-1} .

Let $\mathcal{A} = \bigcup_{t=0}^{\infty} \mathcal{A}_t$. Observe that actually $\mathcal{A} = \mathcal{A}_n$, i.e., the process becomes stationary after at most n steps.

Lemma 5.11 *W.h.p. the outcome of the above process satisfies $|\mathcal{A}| = |V \setminus U| - o(n)$.*

Proof. Let \mathcal{A}_t be the set constructed after t steps of the above process, with $\mathcal{A}_0 = \mathcal{C}$ and $\mathcal{A}_{-1} = \emptyset$. Let \mathcal{H}_t be the history of the process up to time t . Let $v \in V \setminus (U \cup \mathcal{A}_t)$ be a vertex, and let $e_v \in \mathcal{M}$ be the random edge supported by v . The only conditioning that \mathcal{H}_t imposes on e_v is that e_v has at least one vertex $w \neq v$ that does not lie in \mathcal{A}_{t-1} . Hence,

$$\mathbb{P}[v \notin \mathcal{A}_{t+1} | \mathcal{H}_t] = \mathbb{P}[e_v \setminus \{v\} \not\subset \mathcal{A}_t | \mathcal{H}_t] \leq (k-1) \cdot \frac{|V \setminus (U \cup \mathcal{A}_t)|}{|V \setminus (U \cup \mathcal{A}_{t-1})|}. \quad (38)$$

To analyze the quantity on the right, let $a_t = |\mathcal{A}_t| / |V \setminus U|$ for $t \geq -1$. Then Lemma 5.10 implies that w.h.p. $a_0 \geq 1 - \lambda^l \exp(-\lambda)$. With this notation, (38) reads

$$\mathbb{E}[1 - a_{t+1} | \mathcal{H}_t] \leq \frac{(k-1)(1-a_t)^2}{1-a_{t-1}}.$$

Furthermore, given \mathcal{H}_t , for all vertices $v \in V \setminus (U \cup \mathcal{A}_t)$ the events $\{v \notin \mathcal{A}_{t+1}\}$ are mutually independent (because each is determined by the edge e_v supported by v). Therefore, the number of $v \in V \setminus (U \cup \mathcal{A}_t)$ such that $v \notin \mathcal{A}_{t+1}$ is stochastically dominated by a binomial distribution with mean $|V \setminus U| \cdot \frac{(k-1)(1-a_t)^2}{1-a_{t-1}}$. By Chernoff bounds, with probability $1 - o(1/n)$ we therefore see that the number of $v \in V \setminus (U \cup \mathcal{A}_t)$ such that $v \notin \mathcal{A}_{t+1}$ is $|V \setminus U| \cdot \frac{(k-1)(1-a_t)^2}{1-a_{t-1}} + o(n)$. Hence,

$$\mathbb{P}\left[a_{t+1} < 1 - \frac{(k-1)(1-a_t)^2}{1-a_{t-1}} + o(1) | \mathcal{H}_t\right] = o(1/n), \quad (39)$$

and thus the above holds for all $t \geq 1$ w.h.p.

Now, consider the (deterministic) recurrence

$$\alpha_0 = \lambda^l \exp(-\lambda), \quad \alpha_{t+1} = 1 - \frac{(k-1)(1-\alpha_t)^2}{1-\alpha_{t-1}}.$$

It is straightforward to verify that $\lim_{t \rightarrow \infty} \alpha_t = 1$. Therefore, (39) implies that w.h.p.

$$\lim_{t \rightarrow \infty} |V \setminus (U \cup \mathcal{A}_t)|/n = 0,$$

and thus $|V \setminus (U \cup \mathcal{A})| = o(n)$ w.h.p. \square

Proof of Proposition 5.4. We are left to show that w.h.p. all vertices in \mathcal{A} are n/k^3 -rigid. We start by proving that w.h.p. all vertices in \mathcal{C} are n/k^3 -rigid. Suppose that there is another 2-coloring τ of \mathcal{C} such that the set

$$\Delta = \{v \in \mathcal{C} : \sigma(v) \neq \tau(v)\}$$

has size $0 < |\Delta| < n/k^3$. By the construction of \mathcal{C} , each vertex $v \in \Delta$ supports at least 3 edges that consist of vertices in \mathcal{C} only. As these edges are bichromatic under τ , each of them must contain a second vertex in Δ . Hence, there are at least $3|\Delta|$ edges that are supported by a vertex in Δ (under σ) and that contain a second vertex in Δ . But Lemma 5.9 shows that w.h.p. there is no such set Δ of size $0 < |\Delta| < n/k^3$. This shows that all vertices in \mathcal{C} are n/k^3 -rigid w.h.p.

Furthermore, the construction of \mathcal{A} ensures that any 2-coloring τ of H_1 such that $\tau(v) \neq \sigma(v)$ for some $v \in \mathcal{A}$ is indeed such that $\tau(w) \neq \sigma(w)$ for some $w \in \mathcal{C}$. This shows that any $v \in \mathcal{A}$ is n/k^3 -rigid w.h.p., because any $w \in \mathcal{C}$ is. \square

6 A closer look at the internal entropy: proof of Corollary 1.4

6.1 Outline

Throughout this section, we let $f_0(n)$ denote a function such that $f_0(n) = o(n)$ as $n \rightarrow \infty$. Let $\sigma \in \{0, 1\}^n$ be such that $||\sigma^{-1}(0)| - |\sigma^{-1}(1)|| \leq f_0(n)$. In addition, let σ_0 be an equitable 2-coloring. To prove Corollary 1.4 we need to prove that the size $|\mathcal{C}(\sigma)|$ of the local cluster in the planted model $H_k(n, m, \sigma)$ is tightly concentrated. To accomplish that, we need to study the set U from Section 5. That is, $U \subset V$ is constructed as follows.

1. Initially, let U consist of all vertices that do not support any edges.
2. While there is a vertex $v \notin U$ that does not support an edge that does not contain a vertex from U , add v to U .

As a first step, we are going to show that $|U|$ is tightly concentrated. More precisely, in Section 6.2 we will prove the following.

Proposition 6.1 *For any two functions $f_0(n) = o(n)$, $f_1(n) = o(n)$ there is a function $f_2(n) = o(n)$ such that*

$$\mathbb{P} [||U| - \mathbb{E}_{H_k(n, m, \sigma_0)}|U|| > f_2(n)] \leq \exp(-f_1(n)).$$

We also need the following simple expansion properties.

Lemma 6.2 *W.h.p. both $H_k(n, m)$ and $H_k(n, m, \sigma)$ have the following property.*

$$\text{For any set } S \subset V \text{ of size } |S| \leq 2^{-k^2}n \text{ the number of edges } e \text{ that contain at least two vertices from } S \text{ is bounded by } 1.01|S|. \quad (40)$$

Furthermore, with probability $1 - \exp(-\Omega(n))$, $H_k(n, m, \sigma)$ has the following property.

$$\text{For any set } S \subset V \text{ of size } 2^{-k^2}n < |S| \leq n/k^3 \text{ the number of critical edges } e \text{ that contain at least two vertices from } S \text{ is bounded by } 1.01|S|. \quad (41)$$

Proof. This follows from a simple first moment argument similar to the one in the proof of Lemma 5.9. \square

Using Proposition 6.1 and Lemma 6.2, we will derive the following in Section 6.3.

Proposition 6.3 Let $\nu_k(n, m) = \mathbb{E}_{H_k(n, m, \sigma_0)} \ln \mathcal{C}(\sigma_0)$. For any $f_0(n), f_1(n) = o(n)$ there is a function $f_3(n) = o(n)$ such that

$$\mathbb{P}_{H_k(n, m, \sigma)} [|\nu_k(n, m) - \ln \mathcal{C}(\sigma)| > f_3(n) \text{ and (40) holds}] \leq \exp(-f_1(n)). \quad (42)$$

Furthermore, for any $1 \leq j \leq n^{2/3}$ we have

$$0 \leq \nu_k(n, m) - \nu_k(n, m + j) = o(n^{3/4}). \quad (43)$$

Proof of Corollary 1.4. To begin, let us fix a small $\varepsilon > 0$. Our first goal is to show that there exists a density $r = r(n)$ such that

$$\frac{1}{n} \mathbb{E}_{H_k(n, m, \sigma)} \ln |\mathcal{C}(\sigma)| \sim \frac{1}{n} \ln \mathbb{E}_{H_k(n, m)} Z - \varepsilon. \quad (44)$$

To prove (44), it is easier to work with the random hypergraph $H_k(n, p, \sigma)$ in which each $e \subset V$ of size k that is bicolored under σ is inserted with probability p independently. Then for any fixed n , the function

$$F_n(p) = \frac{1}{n} \mathbb{E}_{\sigma, H_k(n, p, \sigma)} \ln |\mathcal{C}(\sigma)|$$

is a polynomial in p . Furthermore, it is clear that $F_n(p) \rightarrow \ln 2$ as $p \rightarrow 0$, and $F_n(p) \rightarrow o(1)$ as $p \rightarrow 1$. For any p we let $\rho(p) \geq 0$ be such that the expected number of edges in $H_k(n, p, \sigma)$ equals $\rho(p)n$. Then by the mean value theorem, there exists p such that $F_n(p) \sim \frac{1}{n} \ln \mathbb{E}_{H_k(n, \lceil \rho(p)n \rceil)} Z - \varepsilon$. Since the acutal number of edges of $H_k(n, p, \sigma)$ is binomially distributed and therefore tightly concentrated about $\rho(p)n$, the ‘continuity property’ (43) ensures that

$$\frac{1}{n} \mathbb{E}_{H_k(n, \lceil \rho(p)n \rceil, \sigma)} \ln |\mathcal{C}(\sigma)| \sim F_n(p) \sim \frac{1}{n} \ln \mathbb{E}_{H_k(n, \lceil \rho(p)n \rceil)} Z - \varepsilon.$$

Setting $r_\varepsilon(n) = \rho(p(n))$, we obtain (44).

For this density $r = r_\varepsilon(n)$ there exists a function $f_1(n) = o(n)$ such that

$$\ln(g_{k, n, m}[\mathcal{B}]) \leq \ln(p_{k, n, m}[\mathcal{B}]) + f_1(n) \quad \text{for any event } \mathcal{B} \neq \emptyset. \quad (45)$$

Let $f_0(n)$ be such that with probability $1 - \exp(-2f_1(n))$, a random $\sigma \in \{0, 1\}^n$ satisfies $\|\sigma^{-1}(0) - \sigma^{-1}(1)\| \leq f_0(n)$. Combining Lemma 6.2, Proposition 6.3, (44) and (45), we see that for these densities $r_\varepsilon(n)$, w.h.p. a random pair (H, σ) chosen from the Gibbs distribution is such that

$$\frac{1}{n} \ln |\mathcal{C}(\sigma)| \geq \frac{1}{n} \ln \mathbb{E}[Z] - \varepsilon \geq \frac{1}{n} \ln Z(H) - 2\varepsilon. \quad (46)$$

Since (46) holds w.h.p. for any fixed $\varepsilon > 0$, there exist sequences $\varepsilon(n) \rightarrow 0$, $r(n)$ as desired. \square

6.2 Proof of Proposition 6.1

We are going to trace the process for the construction of the set U via the method of differential equations [21]. To obtain sufficient concentration from this approach, we will have to modify the process slightly. The modified process will yield a subset $U_* \subset U$, whose size is tightly concentrated. We will then see how U_* can be enhanced to a superset $U^* \supset U$, whose size does not exceed the size of U_* significantly with a very high probability.

Our construction of U_* comes with a parameter $\omega \geq \omega_0$, where ω_0 denotes a large constant (later we will let $\omega \rightarrow \infty$ slowly as $n \rightarrow \infty$). To construct U_* , we consider a similar process as in the proof of Lemma 5.7, but we only run this process on the set V' of vertices that support at most ω clauses. In each step, any vertex $w \in V'$ is either alive, dead, or neutral. Initially, all vertices in V' that do not support a clause are alive, and all others are neutral. The process stops once there is no alive vertex left. In each step, an alive vertex v is chosen randomly. Let d_v be the number edges e_1, \dots, e_{d_v} supported by neutral vertices in which v occurs.

Case 1: $d_v \leq \omega$. All of e_1, \dots, e_{d_v} are deleted from the hypergraph.

Case 2: $d_v > \omega$. In this case ω edges amongst e_1, \dots, e_{d_v} are chosen randomly and are deleted from the hypergraph. Moreover, the remaining $d_v - \omega$ edges are *changed* as follows. Suppose that the deleted edges are e_1, \dots, e_ω . Then v is replaced in each edge $e \in \{e_{\omega+1}, \dots, e_{d_v}\}$ independently by a random vertex $w \neq v$ with $\sigma(w) = \sigma(v)$ that is not dead and that does not belong to e already; if there is no such vertex w left, the process stops.

Finally, all neutral vertices that do not support an edge anymore (after the edge deletions described above) are declared alive, and v is declared dead. Let T be the stopping time of the process, and let U_* be the set of dead vertices upon termination. Then $|U_*| = T$.

The difference between the above process and the actual construction of U is that the latter runs on the entire set V (not just V') and that it *always* removes the e_1, \dots, e_{d_v} . Therefore, $U_* \subset U$.

To trace the construction of U_* , we need to define a few random variables. For each $1 \leq s \leq \omega$ and each $1 \leq l \leq s$ let $X_t(s, l)$ denote the number of neutral vertices that support s vertices in total, out of which l do not contain a vertex that has died by the end of step t . In addition, let A_t signify the number of alive vertices. Let $(\mathcal{F}_t)_{t \geq 0}$ be the filtration generated by the random variables $X_t(s, l)$ and A_t .

Let \mathcal{D}_s be the number of vertices that support precisely s edges ($s \geq 0$). Moreover, let $\mathcal{D}_{>\omega}$ be the number of vertices that support more than ω edges.

Lemma 6.4 *We have*

$$\mathbb{P}[\mathcal{D}_{>\omega} > \exp(-\omega)n] \leq \exp\left[-\frac{n}{2\exp(2\omega)r}\right].$$

Furthermore, for any $0 \leq s \leq \omega$ we have

$$\mathbb{P}[|\mathcal{D}_s - \mathbb{E}\mathcal{D}_s| > \exp(-\omega^2)n] \leq \exp\left[-\frac{n}{2\exp(2\omega^2)r}\right].$$

Proof. For each vertex v the number $s(v)$ of edges supported by v has a binomial distribution with mean $\lambda = kr/(2^{k-1} - 1)$. Assuming that $(k$ and thus) λ is sufficiently large, and choosing ω_0 big enough, we see from Chernoff bounds that $\mathbb{P}[s(v) > \omega] \leq \exp(-8\omega)$. Hence, $\mathbb{E}\mathcal{D}_{>\omega} \leq n \exp(-8\omega)$. Furthermore, $\mathcal{D}_{>\omega}$ satisfies a Lipschitz condition: adding or removing a single edge can alter the value of $\mathcal{D}_{>\omega}$ by at most one. Therefore, the first assertion follows from Azuma's inequality. Similarly, adding or removing a single edge can change the value of \mathcal{D}_s by at most one, and thus Azuma's inequality also implies the second claim. \square

Lemma 6.5 *For any $1 \leq t < \min\{T, n/k^2\}$ we have*

$$\mathbb{E}[X_{t+1}(s, l) | \mathcal{F}_t] = X_t(s, l) \left(1 - \frac{l(k-1)}{n-t}\right) + X_t(s, l+1) \cdot \frac{(l+1)(k-1)}{n-t} + o_\omega(1) \quad (47)$$

$$\mathbb{E}[A_{t+1} | \mathcal{F}_t] = \frac{k-1}{n-t} \sum_{s=1}^{\omega} X_t(s, 1) + o_\omega(1) \quad (48)$$

Furthermore,

$$|X_{t+1}(s, \lambda) - X_t(s, \lambda)| \leq \omega, \quad |A_{t+1} - A_t| \leq \omega \quad (49)$$

with certainty.

Proof. This is a standard argument for a differential equations analysis, based on the following observation ('method of deferred decisions'): given the history \mathcal{F}_t of the process up to time t , for each neutral vertex w each remaining edge e supported by w is conditioned *only* to the effect that e does not contain a vertex that has died by time t . Thus, the alive vertex v chosen at time $t+1$ has a probability of $1 - (1 - 1/(n-t))^{k-1} \sim$

$(k-1)/(n-t)$ of occurring in each edge supported by a neutral vertex, and these events are independent for all such edges. Furthermore, since each neutral vertex only supports $\leq \omega$ edges (as we confine ourselves to the set V'), the probability that v occurs in two such edges is $o(1)$.

This means that given \mathcal{F}_t the expected number of vertices that support s edges in total, out of which l are left after time t , and which support an edge in which v occurs, equals

$$X_t(s, l) \frac{l(k-1)}{n-t} + o(1). \quad (50)$$

Furthermore, the given \mathcal{F}_t expected number of vertices that support s edges in total with $l+1$ left after time t amongst which precisely one contains v , is

$$X_t(s, l+1) \cdot \frac{(l+1)(k-1)}{n-t} + o(1). \quad (51)$$

To obtain (47) from this, we need to take into account the ‘exceptional’ case 2 of the process. But since the expected number of occurrences of v given \mathcal{F}_t is bounded by $\lambda n/(n-t) \leq 2\lambda$, and since this number is binomially distributed, the probability that v occurs in more than ω edges supported by neutral vertices is bounded by $\exp(-\omega)$. This estimate in combination with (50) and (51) yields (47). Equation (48) follows from a similar argument, and (49) is immediate from the construction. \square

Corollary 6.6 *There exists a number $0 < \mu = \mu(k, r) \leq 2n \exp(-\lambda)$ and a function $\delta_\omega = o_\omega(1)$ such that*

$$\mathbb{P}[|T/n - \mu| \leq \delta_\omega] \geq 1 - \exp\left[-\frac{n}{\exp(\omega^3)}\right]. \quad (52)$$

Proof. Lemma 6.4 and Lemma 6.5 verify the assumptions of [21, Theorem 5.1] for times $t \leq n/k^2$. Furthermore, Proposition 5.3 shows that $T \leq 2n \exp(-\lambda) < n/k^2$ w.h.p. Therefore, we can apply [21, Theorem 5.1] to obtain (52). \square

Remark 6.7 *The ‘method of differential equations’ [21, Theorem 5.1] actually shows that the random variables $X_t(s, l)$ closely trace a system of ordinary differential equations. From these the number μ in Corollary 6.6 could, in principle, be worked out precisely for any given k, r . However, for our purposes it is not important to know μ precisely. In the proof of Corollary 6.6 it is important to use the differential equations approach as in [21] to ensure sufficient concentration.*

Since $|U_*| = T$ and $U_* \subset U$, Corollary 6.6 provides a lower bound on the size of U . As a next step, we will derive an (asymptotically) matching upper bound.

Lemma 6.8 *There is a function $\delta_\omega = o_\omega(1)$ such that*

$$\mathbb{P}[|U \setminus U_*| > \delta_\omega n] \leq 3 \exp\left[-\frac{n}{\exp(\omega^3)}\right].$$

Proof. We use a similar argument as in the proof of Lemma 5.7. Namely, having constructed U_* , we commence a second process. Again, in the course of this process vertices can be alive, dead, or neutral. Initially, all vertices in U_* are dead. Furthermore, a vertex $v \notin U_*$ is declared alive if either $v \in V \setminus V'$ (i.e., v supports more than ω edges), or v supports an edge that contains a vertex that occurs in more than ω edges supported by other vertices. All other vertices are neutral. From this initial state, the process proceeds just like the construction of the set U . Namely, in each step an alive vertex v is chosen, unless there is none left, in which case the process stops. Then, all vertices w such that each edge e supported by w contains

either v or a dead vertex are declared alive, and v dies. Clearly, the set of dead vertices of this process upon termination contains U .

By Corollary 6.6 we may assume that $|U_*| \leq 2n \exp(-\lambda)$. Standard arguments (similar to the proof of Lemma 5.6) show that with probability $\geq 1 - \exp\left[-\frac{n}{\exp(\omega^3)}\right]$ the total number of vertices that are alive initially is $o_\omega(1)n$. Furthermore, using stochastic dominance as in the proof of Lemma 5.7, one can show that the above process will terminate after only $o_\omega(1)n$ steps with probability $\geq 1 - \exp\left[-\frac{n}{\exp(\omega^3)}\right]$. \square

Proof of Proposition 6.1. Corollary 6.6 and Lemma 6.8 show that there exist $\mu > 0$ and for any $\omega \geq \omega_0$ some $\delta_\omega > 0$ such that

$$\mathbb{P}[|U| - \mu n > \delta_\omega n] \leq \exp(-n/\exp(\omega^4)), \quad (53)$$

where $\delta_\omega \rightarrow 0$ as $\omega \rightarrow \infty$. We may assume that the given function $f_1(n) = o(n)$ satisfies $f_1(n) \geq \sqrt{n}$, and that $\delta_\omega \geq 1/\omega$. Given such a $f_1(n)$, we can choose a slowly growing function $\omega = \omega(n) \leq \ln \ln n$ such that $\exp(-n/\exp(\omega^4)) \leq \exp(-f_1(n))$. Then (53) implies that $|\mathbb{E}|U| - \mu n| \leq 2\delta_\omega n$. Thus, setting $f_2(n) = 2\delta_{\omega(n)}$ and invoking (53) once more completes the proof. \square

6.3 Proof of Proposition 6.3

To prove Proposition 6.3, it will be easier to work with a slightly different distribution over the hypergraphs that for which σ is a 2-coloring. Namely, $H_k(n, p, \sigma)$ denote a random hypergraph obtained by including each possible edge that is 2-colored under σ with probability p independently. *Throughout this section, we fix p so that the expected number of edges of $H_k(n, p, \sigma)$ is equal to m .* Due to our assumptions on σ , this means that $p \sim m/((1 - 2^{1-k})\binom{n}{k})$.

Lemma 6.9 *For any event E , $\mathbb{P}[H_k(n, m, \sigma) \in E] \leq O(\sqrt{n})\mathbb{P}[H_k(n, p, \sigma) \in E]$.*

Proof. In $H_k(n, p)$ the total number of edges has a binomial distribution with mean m . Therefore, the probability that $H_k(n, p)$ has exactly m edges is $\Omega(1/\sqrt{m}) = \Omega(1/\sqrt{n})$. Furthermore, given that its total number of edges is m , $H_k(n, p)$ is uniformly distributed over all such hypergraphs for which σ is a 2-coloring. \square

The argument for the proof of Proposition 6.3 basically is as follows. We will see that (essentially) all vertices in $V \setminus U$ are rigid, and thus the entropy of the local cluster stems solely from variations of the colors in U . Furthermore, the hypergraph induced on U by the edges do not already contain two vertices from $V \setminus U$ with different colors is sub-critical, i.e., it decomposes into small (at most $\ln n$ but mostly constant-sized) connected components. Now, for each ‘type’ (i.e., isomorphism class) of component the number of occurrences of this type is tightly concentrated (similarly as in a subcritical random graph). This implies concentration of the total number of colorings on $V \setminus U$ because the total number is simply the sum of the numbers of colorings of the components. Let us now carry out the details.

Consider the following way to construct a set $\mathcal{C} \subset V$ of vertices of $H_k(n, m, \sigma)$ (cf. Section 5.3.). Let $l = 10$.

- C1.** Initially, let \mathcal{C} contain all $v \in V$ that support at least $l/2$ edges.
- C2.** While there is $v \in \mathcal{C}$ that supports $< l/2$ edges consisting of vertices of \mathcal{C} only, remove v from \mathcal{C} .
- C3.** While there is a vertex $v \in V \setminus \mathcal{C}$ that supports an edge e such that $e \setminus \{v\} \subset \mathcal{C}$, add v to \mathcal{C} .

Proposition 6.10 *For any function $g_1(n) = o(n)$ there is a function $g_2(n) = o(n)$ such that with probability $\geq 1 - \exp(-g_1(n))$ the set \mathcal{C} has the following properties.*

1. $|V \setminus \mathcal{C}| = |U| + g_2(n)$.
2. Either (40) is violated, or any 2-coloring $\tau \in \mathcal{C}(\sigma)$ satisfies $\tau(v) = \sigma(v)$ for all $v \in \mathcal{C}$.

Proof. The same arguments as in Section 5.3 apply. \square

For a set $C \subset V$ let $\mathcal{A}(C)$ be the set of all $e \subset V$, $|e| = k$, that have neither of the following two properties.

1. $e \subset C$.
2. There is a color $i \in \{0, 1\}$ such that $|e \cap \sigma^{-1}(i)| = 1$ and $|e \cap \sigma^{-1}(1 - i) \cap C| = k - 1$. (In other words, e is critical with respect to σ and has $k - 1$ vertices, not including the supporting one, in C .)

The reason why it is easier, for the present context, to work with the $H_k(n, p, \sigma)$ model is the following simple observation. If we condition on the outcome $\mathcal{C} \subset V$ of the process **C1–C3**, each $e \in \mathcal{A}(C)$ is present as an edge in $H_k(n, p, \sigma)$ with probability p independently. That is, the distribution of $H_k(n, p, \sigma)$ *outside* the ‘core’ \mathcal{C} can be captured very easily.

Given the outcome \mathcal{C} of the process **C1–C3**, let $\bar{H}_k(n, p, \sigma, \mathcal{C})$ denote the random hypergraph on $V \setminus \mathcal{C}$ in which we include the set $e \setminus \mathcal{C}$ for each edge e of $H_k(n, m, \sigma)$ such that $|e \setminus \mathcal{C}| \geq 2$.

Lemma 6.11 *Suppose $|V \setminus \mathcal{C}| \leq 3 \exp(-\lambda)n$. W.h.p. all connected components of $\bar{H}_k(n, p, \sigma, \mathcal{C})$ have size $O(\ln n)$. Furthermore, for any $\omega = \omega(n) \rightarrow \infty$ the expected number of vertices of $\bar{H}_k(n, p, \sigma, \mathcal{C})$ that belong to components of size at least ω is bounded by $\exp(-\Omega(\omega))n$.*

Proof. Given the above observation, the assertion is a direct consequence of the result on the ‘giant component’ phase transition in random non-uniform hypergraphs from [20]. \square

Let \mathcal{T} be the set of all equivalence classes with respect to isomorphism of hypergraphs with edges of size $\leq k$. An *isolated copy* of $T \in \mathcal{T}$ in $\bar{H}_k(n, p, \sigma, \mathcal{C})$ is a subset $S \subset V \setminus \mathcal{C}$ such that S is a component of $\bar{H}_k(n, p, \sigma, \mathcal{C})$ and such that the sub-hypergraph induced on S is isomorphic to T . Let $Y_{T, \mathcal{C}}$ signify the number of isolated copies of $T \in \mathcal{T}$ in $\bar{H}_k(n, p, \sigma, \mathcal{C})$ (given the set \mathcal{C}).

Lemma 6.12 *For any $T \in \mathcal{T}$ and any $d > 0$ we have*

$$\mathbb{P}[|Y_{T, \mathcal{C}} - \mathbb{E}Y_{T, \mathcal{C}}| > d] \leq \exp\left(-\frac{d^2}{16k^2m}\right). \quad (54)$$

Furthermore, if $|V \setminus \mathcal{C}| \leq 3 \exp(-\lambda)n$, then for any $\omega = \omega(n) \rightarrow \infty$ we have

$$\sum_{T \in \mathcal{T}: |V(T)| \leq \omega} |V(T)| \cdot \mathbb{E}[Y_{T, \mathcal{C}}] \geq (1 - \exp(-\Omega(\omega)))n.$$

Proof. The random variable Y_T satisfies a Lipschitz condition: either adding or removing an edge to/from $\bar{H}_k(n, p, \sigma, \mathcal{C})$ can change Y_T by at most k . Therefore, the first assertion follows from Azuma’s inequality. The second one is an immediate consequence of (54) and Lemma 6.11. \square

We will now drop the conditioning upon the outcome of the process **C1–C3**. That is, we let $\bar{H}_k(n, p, \sigma)$ be the random hypergraph obtained by first constructing \mathcal{C} in $H_k(n, p, \sigma)$ and then performing the construction of $H_k(n, p, \sigma, \mathcal{C})$. For each $T \in \mathcal{T}$ let Y_T be the number of isolated copies of T in $\bar{H}_k(n, p, \sigma)$.

Corollary 6.13 *For any function $g_1(n) = o(n)$ and any $\omega = \omega(n) \rightarrow \infty$ there exists $g_2(n) = o(n)$ such that the following is true. For each $T \in \mathcal{T}$ there is a number $y_T = y_T(k, r) \geq 0$ such that with probability $\geq 1 - \exp(-g_1(n))$ either (40) is violated or the following is true.*

1. All but $g_2(n)$ vertices of $\bar{H}_k(n, p, \sigma)$ belong to a component on $\leq \omega$ vertices.
2. We have $\sum_{T \in \mathcal{T}} |V(T)| \cdot |Y_T - y_T n| \leq 2g_2(n)$.

Proof. This is immediate from Proposition 6.10 (which, crucially, shows that $|\mathcal{C}|$ is tightly concentrated) and Lemma 6.12. \square

For each $T \in \mathcal{T}$ let z_T denote the number of 2-colorings of T . Furthermore, let $Z(\bar{H}_k(n, p, \sigma))$ denote the number of 2-colorings of $\bar{H}_k(n, p, \sigma)$.

Corollary 6.14 *For any function $g_1(n) = o(n)$ there exists $g_2(n) = o(1)$ such that with probability $\geq 1 - \exp(-g_1(n))$ either (40) is violated or we have*

$$\left| \frac{1}{n} \ln Z(\bar{H}_k(n, p, \sigma)) - \sum_{T \in \mathcal{T}} y_T z_T \right| \leq g_2(n).$$

Proof of Proposition 6.3. Let us first deal with the random hypergraph $H_k(n, p, \sigma)$. Suppose that $|\mathcal{C}| \geq n(1 - 2\exp(-\lambda))$. Then any 2-coloring τ of $\bar{H}_k(n, p, \sigma, \mathcal{C})$ yields an element of the local cluster $\mathcal{C}_{H_k(n, p, \sigma, \mathcal{C})}(\sigma)$ of $H_k(n, p, \sigma, \mathcal{C})$ by letting $\tau(v) = \sigma(v)$ for all $v \in \mathcal{C}$. Therefore, Proposition 6.10 and Corollary 6.14 imply that

$$\mathbb{P}_{H_k(n, p, \sigma)} \left[\frac{1}{n} \ln |\mathcal{C}(\sigma)| \geq \sum_{T \in \mathcal{T}} y_T z_T - o(1) \right] \geq 1 - \exp(-f_1(n)). \quad (55)$$

Conversely, consider a coloring τ of $H_k(n, p, \sigma)$. By Proposition 6.10, either (40) is violated or $\tau(v) = \sigma(v)$ for all $v \in \mathcal{C}$. Assume that the latter is true. Then τ induces a 2-coloring of $\bar{H}_k(n, p, \sigma, \mathcal{C})$.

$$\frac{1}{n} \ln \mathcal{C}_{H_k(n, p, \sigma)}(\sigma) \leq \sum_{T \in \mathcal{T}} y_T z_T + o(1). \quad (56)$$

Hence,

$$\mathbb{P}_{H_k(n, p, \sigma)} \left[\text{either (40) is violated or } \frac{1}{n} \ln |\mathcal{C}(\sigma)| \leq \sum_{T \in \mathcal{T}} y_T z_T + o(1) \right] \geq 1 - \exp(-f_1(n)). \quad (57)$$

Combining (55) and (57) with Lemma 6.9, we obtain (42).

Finally, to obtain (43), observe that adding a further of $n^{2/3}$ edges to $H_k(n, m, \sigma)$ will simply can just connect at most $n^{2/3}$ components of $\bar{H}_k(n, p, \sigma, \mathcal{C})$ with the set \mathcal{C} . Lemma 6.11 shows that w.h.p. all of these components have size $\leq n^{0.01}$. Hence, w.h.p. the total reduction in the number of 2-colorings is $\leq n^{2/3+0.01} = o(n^{3/4})$. \square

References

- [1] D. Achlioptas, A. Coja-Oghlan: Algorithmic barriers from phase transitions. Proc. 49th FOCS (2008) 793–802.
- [2] D. Achlioptas, J.H. Kim, M. Krivelevich, P. Tetali: Two-coloring random hypergraphs. Random Structures and Algorithms **18** (2002), 249–259.
- [3] D. Achlioptas, C. Moore: Random k -SAT: two moments suffice to cross a sharp threshold. SIAM Journal on Computing **36** (2006) 740–762.

- [4] D. Achlioptas, A. Naor: The two possible values of the chromatic number of a random graph. *Annals of Mathematics* **162** (2005), 1333–1349.
- [5] D. Achlioptas, Y. Peres: The threshold for random k -SAT is $2^k \ln 2 - O(k)$. *Journal of the AMS* **17** (2004) 947–973.
- [6] D. Achlioptas, F. Ricci-Tersenghi: On the solution space geometry of random constraint satisfaction problems. *Proc. 38th STOC* (2006) 130–139.
- [7] M. Bayati, D. Gamarnik, P. Tetali: Combinatorial approach to the interpolation method and scaling limits in sparse random graphs. *Proc. 42nd STOC* (2010) 105–114.
- [8] A. Coja-Oghlan, A. Frieze: Random k -SAT: the limiting probability for satisfiability for moderately growing k . *Electronic Journal of Combinatorics* **15** (2008) N2.
- [9] L. Dall’Asta, A. Ramezanpour, R. Zecchina: Entropy landscape and non-Gibbs solutions in constraint satisfaction problems. *Phys. Rev. E* **77**, 031118 (2008).
- [10] O. Dubois, J. Mandler: The 3-XORSAT threshold. *Proc. 43rd FOCS* (2002) 769–778.
- [11] E. Friedgut: Sharp thresholds of graph properties, and the k -SAT problem. *Journal of the AMS* **12** (1999) 1017–1054.
- [12] E. Friedgut: Hunting for sharp thresholds. *Random Struct. Algorithms* **26** (2005) 37–51
- [13] A. Frieze, N. Wormald: Random k -Sat: a tight threshold for moderately growing k . *Combinatorica* **25** (2005) 297–305.
- [14] S. Janson, T. Łuczak, A. Ruciński: *Random Graphs*, Wiley 2000.
- [15] W. Kauzmann: The nature of the glassy state and the behavior of liquids at low temperatures. *Chem. Rev.* **43** (1948) 219–256.
- [16] F. Krzakala, A. Montanari, F. Ricci-Tersenghi, G. Semerjian, L. Zdeborová: Gibbs states and the set of solutions of random constraint satisfaction problems. *Proc. National Academy of Sciences* **104** (2007) 10318–10323.
- [17] A. Montanari, R. Restrepo, P. Tetali: Reconstruction and clustering in random constraint satisfaction problems. *arXiv:0904.2751v1* (2009).
- [18] T. Mora, L. Zdeborová: Random subcubes as a toy model for constraint satisfaction problems. *J. Stat. Phys.* **131** (2008) 1121–1138.
- [19] B. Pittel, G. Sorkin: The satisfiability threshold for k -XORSAT. Preprint (2011).
- [20] Schmidt-Pruzan, J., Shamir, E.: Component structure in the evolution of random hypergraphs. *Combinatorica* **5** (1985) 81–94
- [21] N. Wormald: The differential equation method for random graph processes and greedy algorithms. In M. Karoński and H.J. Prömel (eds.): *Lectures on Approximation and randomized algorithms* (1999) 73–155.
- [22] L. Zdeborová, M. Mézard: Constraint satisfaction problems with isolated solutions are hard. *J. Stat. Mech.* P12004 (2008).